



Deelnemerslijst
Data Management Fundamentals
17-19 november 2025

BEDRIJF	NAAM DEELNEMER		FUNCTIE
APG	Shanmukha	Masapu	business-analyst
Aegon Investment Management B.V.	André	Righolt	business-analyst
All Your BI	Ronald	van Gils	data architect
Bank Indonesia	Laras	Pramesti	data architect
European Investment Bank	Roberto	Alvarez	business-analyst
FMO (Nederlandse Financierings-Maatschappij voor Ontwikkelingslanden N.V.)	Stephanie	Allegrini	business-analyst
NAM	Duke	IGBUWE	it-manager

**Evaluation Form
Data Management Fundamentals
November 17-19, 2025**

Name: _____ Company: _____

What would be your overall grade on a 1 (worst) to 10 (best) scale (please tick):

1. The workshop:

1	2	3	4	5	6	7	8	9	10
<input type="radio"/>									

2. The speaker:

1	2	3	4	5	6	7	8	9	10
<input type="radio"/>									

- Content:

- Presentation skills:

3. Did the programme live up to your **expectations?** **Yes** **Partially** **No, because**

4. Would you recommend this workshop to colleagues or peers? **Yes** **No** (please explain)

5. Which subjects did you miss or were not covered adequately?

6. Which subjects were superfluous or took up too much time in your opinion?

7. What additional comments can you offer?

8. How would you grade the organisation and venue:

1	2	3	4	5	6	7	8	9	10
----------	----------	----------	----------	----------	----------	----------	----------	----------	-----------

Overall organisation of the workshop

Quality classroom, screen and sound

Geographic location / accessibility

9. How did you come here? **Car** **Public transport** **Other** _____

10. Please leave your review/recommendation if you like, which we may post on our website.

Name: _____

Job title: _____ Company: _____

Welcome Adept Events

WHO WE ARE

AdeptEvents

DW&BI SUMMIT

BI-Platform.

RELEASE.

Werner Schoots

Founder Adept Events



BI-Platform.

- Launched in 2008 as online spin-off from Database Magazine (DB/M)
- Topics: Business Intelligence, Data Warehousing, Analytics, Data Management

News

Job board

Selected Whitepapers

Events

Articles

Blogs

Video interviews

Cases

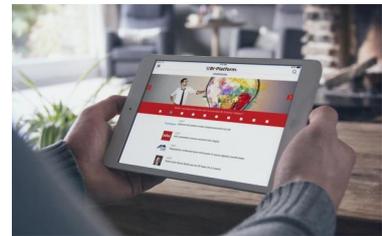
- We welcome your input: redactie@biplatform.nl

  www.biplatform.nl & weekly newsletter

  @BIplatform on & YouTube

  Download the BI-Platform App

 Join our LinkedIn Discussion Group



RELEASE

- Launched in 1996 as Software Development spin-off from Database Magazine
- Topics: Software Engineering – Analysis, Design, Development, Testing and Deployment

News

Job board

Selected Whitepapers

Events

Articles

Blogs

Video interviews

Cases

- We welcome your input: redactie@release.nl

  www.release.nl & weekly newsletter

  @Release_nl on Twitter & YouTube

  Download the Release App

 Join our LinkedIn Discussion Group



SEMINARS

All seminars and workshops are organised twice a year

Alec Sharp	Business-oriented Data Modelling Masterclass Working with Business Processes Masterclass Concept Modelling for Business Analysts The Data-Process Connection (virtual half day session)
Panos Alexopoulos	Knowledge Graphs – pragmatic approach and best practices
Rick van der Lans	Ontwerpen van een Nieuwe Data Architectuur
Mathias Vercauteren	Data Governance Sprint
Nigel Turner	A Data Strategy for Becoming Data Driven Tackling Data Quality Problems (virtual half day session)
Chris Bradley / Winfried Etzel	Data Management Fundamentals
Lawrence Corr	Agile Data Warehouse Design & Dimensional Modeling
Christian Gijssels	Generatieve-AI in Business Analyse Cursus Sparx Enterprise Architect 16
<i>Multiple speakers</i>	<i>Data Warehousing & BI Summit – Yearly conference in March/April</i>

IN-HOUSE

All seminars and workshops can be organized in-company.
With local speakers and international speakers!



Please contact Werner Schoots

☎ +31 (0)172 742680

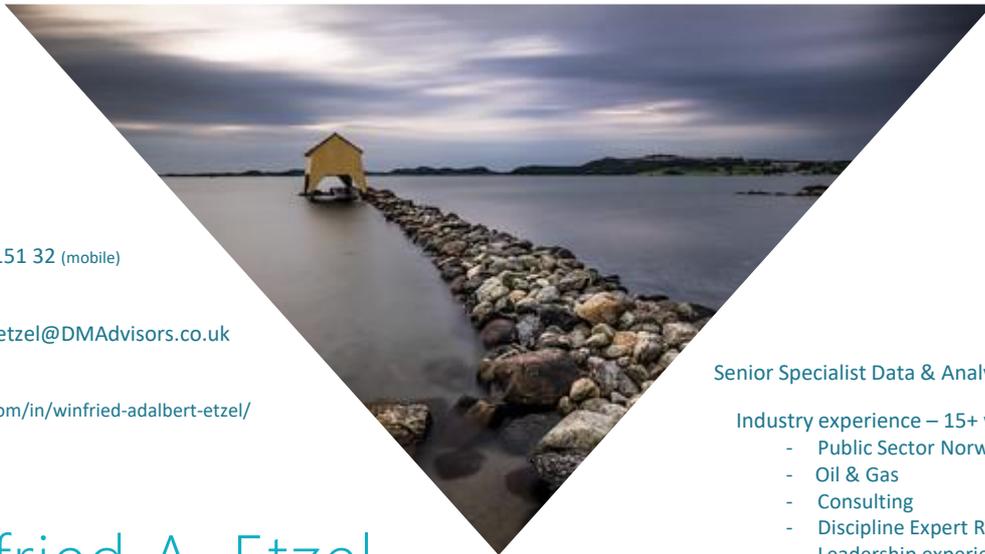
✉ seminars@adeptevents.nl



DMBoK2, CDMP Data Management Fundamentals

WINFRIED A. ETZEL CDMP Master

winfried.etzel@dmadvisors.co.uk



+47 942 151 32 (mobile)

winfried.etzel@DMAAdvisors.co.uk

[LinkedIn.com/in/winfried-adalbert-etzel/](https://www.linkedin.com/in/winfried-adalbert-etzel/)

Senior Specialist Data & Analytics - Equinor

Industry experience – 15+ years

- Public Sector Norway
- Oil & Gas
- Consulting
- Discipline Expert Roles
- Leadership experience in Data

Winfried A. Etzel

INFORMATION MANAGEMENT STRATEGIST



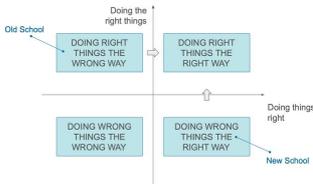
Winfried A. Etzel

Winfried Adalbert Etzel is a key figure in the **Nordic data community**, known for his strong contributions to professional forums and his commitment to promoting Data Governance, Data Strategy, and organizational design for data teams.

Winfried has dedicated his career to emphasizing the importance of information and data management, as well as a strategic and holistic approach to treating information and **data as a valuable asset** for organizations.

Through his work with the **#MetaDAMA podcast** and his voluntary engagement in the Nordic data community, he has been a driving force for innovation and has helped build a strong data community while promoting professional development across the Nordic region.

Winfried is currently working on his book under the working title **«Data Governance in the Wild»**, where he explores how Data Governance needs to change to adapt to distributed landscapes and the need for automation.



Winfried Adalbert Etzel

Home Links About

The Ethics of AI and Data - Human at the Center
 In our latest episode of #MetaDAMA, we discuss data ethics, AI, and the broader societal implications of technology.
 Feb 28 · 15 min

The Role of Data Leadership in the Industrial Sector
 Data leadership in industrial organizations.
 Feb 28 · 15 min

Data Governance—defined for “the Wild”?
 Data Governance is a human-based system by which data assets in a socio-technical system are described, measured and by which...
 Feb 28 · 15 min

META DAMA Podcast
MetaDAMA - Data Management in the Nordics
 Winfried Adalbert Etzel - DAMA Norway

All Episodes

- #M01 - Old School: The Ethics of AI and Data - Human at the Center (2nd)
- #M02 - Old School: The Ethics of AI and Data - Human at the Center (1st)
- #M03 - Old School: The Role of Data Leadership in the Industrial Sector (2nd)
- #M04 - Old School: The Role of Data Leadership in the Industrial Sector (1st)
- #M05 - Old School: The Role of Data Leadership in the Industrial Sector (1st)

About
 This META DAMA Norway's podcast is created as a series for sharing experiences within Data Management, discussing concepts and practical challenges in the field. The goal is to help you get up to speed with your data management journey. The podcast is hosted by Winfried Adalbert Etzel, a professional in the field of Data Management. The podcast is produced by DAMA Norway. The podcast is available on various platforms including Apple Podcasts, Spotify, and Amazon Music. The podcast is also available on the MetaDAMA website. The podcast is a valuable resource for anyone interested in Data Management, Data Governance, and Data Strategy. The podcast is a must-listen for anyone looking to stay up to date on the latest trends in the field. The podcast is a great way to learn from experts in the field. The podcast is a great way to stay up to date on the latest trends in the field. The podcast is a great way to learn from experts in the field. The podcast is a great way to stay up to date on the latest trends in the field.



Training, Mentoring, and Executive Workshops

A DAMA Registered Education Provider (REP)

We offer training courses for practitioners and executives. Custom-built, training & awareness seminars can also be delivered.



Building a business focused Data Strategy - 3-day course showing the components of a business focused data strategy, how to gain input and buy-in, how to develop realistic implementation roadmaps and transition plans, and what it really means to be "data centric".



Data Management Fundamentals – 5-day intermediate course covering all of the disciplines of Information Management as defined in the DAMA Body of Knowledge (DMBoK) & the changes in DMBoK 2.0. Core aspects & methods of the Information disciplines are explored.

Introduction to Information Management – 1-day high level overview course, introducing the major disciplines of Information Management, why it is critical for business today, and the core subjects within the Information Management topic.



Data Modelling Fundamentals – 3-day intermediate course introducing students to data modelling, its purpose, the different types of models and how to construct and read a data model. & how data models should be used for Business improvement & understanding.



Advanced Data Modeling (including Data Modelling for Big Data) – 3-day advanced course for students with existing data modelling experience to understand the human centric aspects and techniques of data modelling to enable them to build quality models that meet business needs. This also covers major data model patterns & common problem solving. The course also show the applicability of Data Modelling in Big Data, NoSQL, and other non-relational environments.

IM Fundamentals & Practitioner Courses – A series of 1-day (foundation) and 2-3-day (practitioner) Classes to provide practitioners a solid background in a specific Information Management topics. The practitioner workshops explore further detail on the implementation aspects of the Information Management discipline & also cover the CDMP specialist exam syllabus, providing preparation for students seeking to take the specialist CDMP exam.

- Data Modelling Foundation (1 day only)
- Data Governance & Stewardship
- Master & Reference Data Management
- Data Quality Management
- Data Warehouse & Business Intelligence
- Data Integration & Interoperability
- Metadata Management



Executive Workshops – ½ day and 1-day executive workshop(s) designed to give non-technical managers a basic understanding of a various Information Management topics and their importance to the organisation.

Professional Certification – 3-day and 5-day workshop "exam cram" designed to help attendees prepare for sitting professional certifications including DAMA CDMP examinations.



Integrated Business Process, Data & Requirements Definition – 5-day intensive class to show students an integrated requirements discovery and definition approach covering business process, different types of requirements modelling, and the critical role of the conceptual data model.



DMBoK₂ and CDMP® Preparation Live and On-Demand Classes



Data Management Fundamentals

<https://www.dataversity.net/dmbok-and-cdmp-preparation-learning-plan/>

Data Governance

<https://training.dataversity.net/learning-paths/dgs0-dmbok-and-cdmp-preparation-data-governance-specialist-learning-plan>

Data Modelling

<https://training.dataversity.net/learning-paths/dms0-dmbok-and-cdmp-preparation-data-modeling-specialist-learning-plan>

Data Quality Management

<https://training.dataversity.net/learning-paths/dqs0-dmbok-and-cdmp-preparation-data-quality-specialist-learning-plan>

Master & Reference Data Management

<https://www.dataversity.net/>

Metadata Management

<https://www.dataversity.net/>

 CDMP Online learning program
 Approved by DAMA-I
 Part of DATAVERSITY training center
 Based on DMBoK CDMP syllabus

**20% discount code for
DAMA UK
"dmadvisors"**



Please Introduce yourself



1. Where you work & for how long
2. Information Management experience
3. Key learning objective(s) from this course (e.g. are you seeking to take the CDMP DMF exam later)
4. One **interesting non-work** fact about you



CDMP Overview



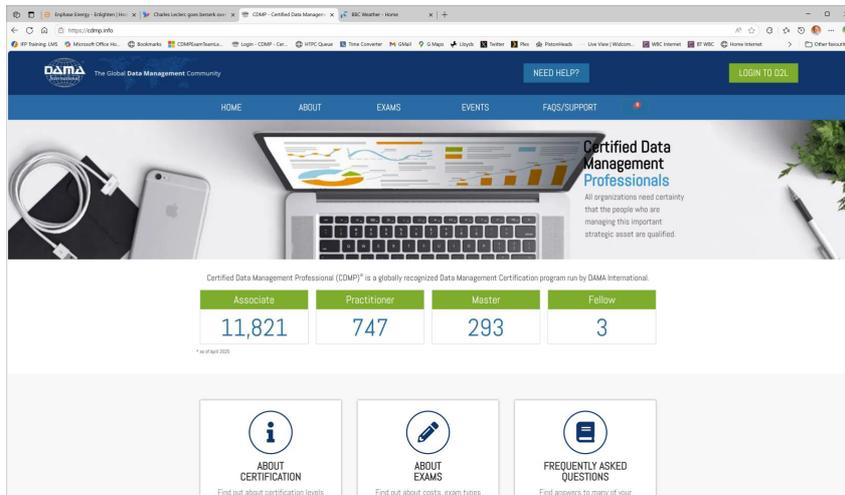
CDMP – Career Path



<https://cdmp.info/>



CDMP Portal: <https://cdmp.info>



Value of the CDMP?

Individual Benefits:

- Professional Development Milestones
- Continuous Professional Growth
- Internationally Recognized Measurement of Accomplishment
- Greater credibility and recognition in the marketplace
- May be a requirement in some areas

Organizational Benefits:

- A benchmark for assessing Data Management practices
- Increased trust in your Data Management team
- Employee development
- Demonstration of commitment to quality standards
- Raise awareness of the importance of Data Management disciplines to the organization



CDMP Associate

- 1 exam – **Data Management Fundamentals**
- Based upon DMBok2 (see breakdown for %)
- Associate pass mark 60%
- **Immediate** award of CDMP Associate for passing DM Fundamentals with 60%+ pass

Examination

- 100 questions
- Multiple choice .. Only one answer per question
- 90 minutes (+20 mins where English is 2nd Language (ESL))
- 60% to pass



CDMP Practitioner

3 exams

- **Data Management Fundamentals** (70% required to pass)
- **2 Specialist exams** (70% required to pass)

Specialist Exams (each)

- 100 questions
- Multiple choice .. Only one answer per question
- 90 minutes (+20 mins where English is 2nd Language)
- 70% to pass

Specialist Exams

- Data Governance
- Data Modelling
- Data Quality
- Metadata
- Data Integration
- DW / BI
- Master & Reference Data
- *Data Architecture Management (on hold for DMBok 3)*



CDMP Master

3 exams

- **Minimum 10 years Data Management experience** (Evidence of experience reviewed by CDMP Fellows)
- **Data Management Fundamentals** (80% required to pass)
- **2 Specialist exams** (80% required to pass)

Specialist Exams (each)

- 100 questions
- Multiple choice .. Only one answer per question
- 90 minutes (+20 mins where English is 2nd Language)
- 80% to pass

Specialist Exams

- Data Governance
- Data Modelling
- Data Quality
- Metadata
- Data Integration
- DW / BI
- Master & Reference Data
- *Data Architecture Management (on hold for DMBok 3)*



CDMP Levels

	Associate	Practitioner	Master	Fellow
DAMA Membership	Central	Central	Central	Central
Practical Experience	Guideline²	2 years – 10 years	10 years <i>minimum</i> ²	Over 25 years ²
	Maximum	N/A ¹	N/A ¹	N/A ¹
Exams	1 DM Fundamentals	3 DM Fundamentals + 2 specialist	3 DM Fundamentals + 2 specialist	Globally recognised & respected thought leadership. Well published. Significant contribution to Data Management Profession. CDMP Master. Contribution to CDMP & DMBOK. By nomination.
Allowable substitutions³	No	No	No	
Pass Mark	60%	70%	80%	
Route	Register & take exam	Register & take exams	Register & take exams. Evidence experience (reviewed by CDMP Fellows)	

1. Individuals are welcome to remain at any level without pressure to go upward. You can decide what level you want to be and if you wish to proceed further.
2. Industry experience is a *guideline* for taking the certification. For Master, a minimum level of experience is required.
3. Substitute exams ceased from 1st August 2020.



CDMP Exam Prices

Item	Member
All CDMP exams (taken at conferences, training events, chapters etc..)	\$300
Online proctoring fee	\$11
CDMP Renewal	\$100
Exam re-take	\$200
<p><i>All exams are taken online only (classroom proctored or online proctor)</i> <i>Online or chapter-led proctoring requires full payment up front</i> <i>Please watch out at various international DAMA-I endorsed conferences for exam proctoring and preparation workshops with Pay If You Pass offers.</i></p>	

<https://cdmp.info/>
 support@cdmp.info
 cdmp@DAMA.org



Question bank

Difficulty Level	Knowledge Tested	% of q's in exam pool
Associate level difficulty	DMBoK2 content and understanding	60%
Practitioner level difficulty	Practical Application of DMBoK2	20%
Master level difficulty	Knowledge expected of a professional with 10+ years practical experience in the topic*	20%



*unlikely to be found in DMBoK2

P / 21

Module breakdown: DM Fundamentals

Data Management Process	Ch 1 + 15, 16, 17	2%
Big Data	Ch 14	2%
Data Architecture & L/C management	Ch 4	6%
Document, Records & Content Management	Ch 9	4-5%
Data Ethics	Ch 2	2%
Data Governance	Ch 3	11%
Data Integration & Interoperability	Ch 8	6-7%
Master & Reference Data Management	Ch 10	11%
Data Modelling & Design	Ch 5	11%
Data Quality	Ch 13	11%
Data Security	Ch 7	6%
Data Storage & Operations	Ch 6	4-5%
Data Warehousing & Business Intelligence	Ch 11	11%
Metadata Management	Ch 12	11%



P / 22

Exam Platform



Gmail EDP Data Governan... WBC Website Twitter TweetDeck Fundamentals

Question 3

Which of these are NOT true of Data Governance??

- DG is a continuous process of data improvement
- A DG initiative should always be led by the IT department
- IT is a key stakeholder in DG
- There are different organization models for DG
- DG is the exercise of authority and control over the management of data assets

Question 4

What are the primary responsibilities of a data steward?

- Identifying data problems & issues
- The manager responsible for writing policies and standards that define the data management prog
- Analyzing data quality
- A business role appointed to take responsibility for the quality and use of their organization's data
- The data analyst who is the subject matter expert (SME) on a set of reference data.

P / 3

Page 1:

1 ✓ 2 ✓ 3 ✓

4 ✓ 5 ✓

Page 2:

6 -- 7 ✓ 8 --

9 -- 10 ✓

Page 3:

11 ✓ 12 -- 13 --

14 -- 15 ✓

Page 4:

16 ✓ 17 ✓ 18 --

Exam Platform

- Answered & Unanswered questions clearly indicated
- At submission, warning of unanswered questions listed
- click on these & answer before confirming submission
- As you proceed, note / flag any question to come back to

But always give an answer NOW

Warnings

You have 23 unanswered questions.

- Question 6
- Question 8
- Question 9
- Question 12
- Question 13
- Question 14
- Question 18
- Question 19



DMF PRA 0:32/23 remaining

Page 1:
1 2 3
4 5
Page 2:
6 7 8
9 10
Page 3:
11 12 13
14 15
Page 4:
16 17 18
19 20
Page 5:
21 22 23
24 25

Question 1 (1 point)
In the conceptual data model an instantiation of a particular business entity is described as:

- Dataset
- Row
- Entity occurrence
- Rule
- Record

Question 2 (1 point)
Which of these statements has the most meaningful relationship label?

- An order is connected with order lines.
- An order is associated with order lines.
- An order line contains orders.
- An order is related to order lines.
- An order is composed of order lines.

Question 3 (1 point)
Which of the following statements about business rules is FALSE?

- Action rules are instructions on what to do when data elements contain certain values.
- Data rules constrain how data relates to other data.
- All business rules must be identified prior to the start of the data modelling process.
- Action rules are difficult to define in a data model.
- Data rules cannot be shown on a data model.

Data Manager: x +

https://cdmp.info

CDMP Exam Team

Associate	Practitioner	Master	Total
4,562	412	110	3

Why Data Management Certification?

Purchase exams

Purchase Data Management Fundamentals exam

Data Management Fundamentals Exam

The Data Management Fundamentals exam is required for all CDMP certification levels: Associate, Practitioner, and Master.

Have questions about how the DAMA-DMBOK2 Revised Edition affects your CDMP certification? Click [here](#).

USD\$311

Purchase options

- STANDARD EXAM
- ESL VERSION**
- PRACTICE EXAM ONLY

The practice exam is free with the exam purchase.

Select version (e.g. English as a Second Language)

After purchase, follow the emailed instructions to register an account on the examination platform



This document will assist you to enrol and take your Data Management Fundamentals exam

This enrolment is for: **Data Management Fundamentals Standard exam**

Before you enrol please read the following guidance and rules:

- The email you use to enrol links all your badges and payments. You should, if possible, use the same email to enrol and pay. We recommend using a personal email, rather than a work email, so you don't lose access to your account in future.
- It is against the Canvas Terms of Service for one person to have two or more accounts. If you make a mistake when you enrol you must NOT try to re-enrol as this could create a second account. Please email support@cdmp.com for assistance.
- This information is provided for you only. This document and the links contained must not be shared with anyone. Changing this information can result in you losing your exam results, your certification or access to the platform. If you are found to be using these enrolment instructions without authorisation, your canvas account and exam results will be deleted.



Data Management Fundamentals (online)

- Create a Canvas account and enrol in the exam**
Click the enrolment link to go to Canvas Catalog to enrol. Once you have enrolled, you will receive 2 emails:
 - to confirm your Canvas Account
 - to confirm your enrolment
 Please follow the instructions in both confirmation emails to complete the enrolment process.
- Log into Canvas and try the Practice exam**
The Practice exam is optional but recommended. It is a short exam that draws 40 random questions from a bank of 200 practice questions. You can take the Practice exam as many times as you wish.

Once you have completed the enrolment process you can access Canvas by logging in here: cdmp.ansistructure.com
- Download Honorlock and take the exam**
Read and accept the exam conditions. You will be prompted to install the Honorlock Chrome Extension. Once you are ready click **Take the Exam**.

If you need assistance during your exam click the chat button for support from Honorlock

The Data Management Fundamentals exam is the foundation exam for all CDMP levels. If you pass this exam with a 70% or more pass mark, you automatically receive an Associate certification. If you achieve a pass mark of 70% it will also count towards Practitioner certification. If you achieve a pass mark of 80% it will also count towards Master certification.

Inside Canvas

Access your exam requirements via the dashboard

Update your notification preferences. We recommend turning off all notifications. You can access these settings via the Account page

Access both the practice exam and the main exam from your course list



Prepare your open book materials

- The open book policy allows only one 'book'. One book means, one of the following:
- DMBK2 in hardcopy, or
 - DMBK2 in digital form, or
 - Hand copy notes, or
 - A digital version of notes
- If you decide to use a digital version of the DMBK2 or your notes, they must be on a separate device, such as a tablet or a separate laptop. You may not use a mobile phone for digital version.

Prepare your exam space

- You will need a quiet desk area, with your face well lit. If windows are behind you, please close the curtains to reduce back light.
- Your desk should be clear of all items except your open book materials
 - Do not have your mobile phone on the desk or in the room
 - Close any doors to the area
 - Cover or remove any writing on walls, such as white boards or pictures with text

Read and watch

- [Exam conditions](#) - you will be asked to acknowledge these prior to taking the exam
- [Academic Integrity policy](#)
- Browse the [CDMP Code](#)
- Watch the [open exam sample video](#) - not completing an adequate room scan can result in the exam attempt being deemed invalid

Checklist

- Photo ID
- Open book materials
- Quiet room / Internet connection
- Prepare exam space
- Latest version Chrome browser
- Webcam
- Mirror to show your desk if you are unable to move your webcam



I have completed the exam

After the exam click the "Finished the Exam" button from the home page of your exam. For important post exam information.

Take a copy of your records

At the end of each calendar month, exam results are sent to CDMP. If for loading to your membership portal. Once this has occurred you will lose access to your exam. Please take a copy of your exam results, for your records.

Badges and Certificates

If you have passed the exam (a score of 80% or more) you will receive an Associate Badge and Certificate. Badges can take up to 5 business days to be issued. Certificates are issued up to 20 days after badges.

Your badge - You will receive a notification email from Badge once your Badge has been issued. Please follow the link in the email to setup your Badge backpack and share your badges. **Your Certificate** - This will be sent to you at the email you used to create your Canvas account.

Need a Retake?

You can do a retake if you would like to try for a higher pass mark. There is no minimum time frame between your first attempt and your re-take. A retake costs USD200 + practice fees.

[Purchase retake](#)

If you passed with over 70% you are eligible to continue your CDMP journey to Practitioner or Master.

Practitioner requirements

70% pass in Data Management Fundamentals exam and 70% pass in 2 specialist exams

Master requirements*

80% pass in Data Management Fundamentals exam and 80% pass in 2 specialist exams

You can book a specialist exam from [here](#)
Get more questions? [Browse the FAQs](#)

Page 4



Exam conditions

By taking this exam:

You are accepting the exam conditions of taking a CDMP exam as listed below.

You accept that any breaches of these conditions will result in penalties such as having your **exam mark removed**.

You are not permitted to

- Leave the room during the exam
- Use your mobile phone
- Communicate with anyone, except the online proctor
- Wear caps or hats. (Cultural headdress is permitted)
- Have any other people in the room
- Wear headphones or a smartwatch
- Copy exam questions or materials
- Search internet for question answers



Home

Badges

Honorlock

New Analytics

Assignments

Rubrics

Marks

Announcements

Pages

Files

Discussions

Quizzes

Modules

BigBlueButton

Collaborations

Outcomes

Syllabus

People

Settings

All exams are open book exams. The open book policy allows only one 'book'.

One book means, one of the following

- DMBok2 in hardcopy, or
- DMBok2 in digital form, or
- Hard copy notes, or
- A digital version of notes

If you decide to use a digital version of the DMBok or your notes, they must be on a separate device, such as a tablet or a separate laptop. You cannot use a mobile phone.

The device cannot be connected to the computer you are using for the exam.

You cannot use a second monitor on the same computer.

You may **not** use a Mobile phone for digital versions.

Exam browser and system conditions

- Close all tabs in your browser prior to commencing your exam.
- You are only permitted to have the Canvas exam page open.
- All applications except Canvas should be closed during your exam.
- You are permitted to use browser translation if required.
- If you require a dictionary translation website, you will need to submit it to CDMP support for approval **prior** to commencing your exam.
- If the Honorlock proctor attempts to pause your exam to contact you, you must not restart the exam until given permission by the proctor.
- You must respond to all requests by the Honorlock proctor.

What happens if I breach the conditions

What happens if you breach the **Exam conditions** as listed above

1. Your exam mark will be removed.
2. Any badges issued will be revoked.
3. You will be issued with 'Breach of exam conditions' notice.

What happens if you breach the **Verification of identity** conditions as listed above

1. Your account **will** be suspended.
2. Any badges issued will be revoked.
3. You will be issued with a suspension notice.

By clicking on 'Take the exam', you acknowledge and accept the above exam conditions

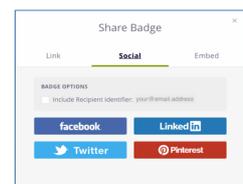
Take the exam



CDMP Badges



- Exam platform awards badges for passing CDMP certifications.
- You can put these badges on your CV, LinkedIn, Twitter, Facebook or even Pinterest.
- Badges are not just an image, Badges carry metadata that can be verified.
- EG, you can send your badge to a prospective employer who can verify that your CDMP certification is authentic and valid.
- Badges are part the global movement called the "OpenBadges" specification.
- Started in 2010 by Mozilla, it has expanded to include millions of Open Badges awarded to hundreds of thousands of recipients & supported by certification players such as The Open Group, PMI, Microsoft and IBM.
- Part of the open badges specification & community are public verification authorities that you can use to independently verify your badges are authentic as required.
- In addition to CDMP badges, accumulate badges for certifications including TOGAF, Prince II, MCSE, PMP & PMI.



DMBoK₂ and CDMP® Preparation

Live and On-Demand Classes



Data Management Fundamentals

<https://www.dataversity.net/dmbok-and-cdmp-preparation-learning-plan/>

Data Governance

<https://training.dataversity.net/learning-paths/dgs0-dmbok-and-cdmp-preparation-data-governance-specialist-learning-plan>

Data Modelling

<https://training.dataversity.net/learning-paths/dms0-dmbok-and-cdmp-preparation-data-modeling-specialist-learning-plan>

Data Quality Management

<https://training.dataversity.net/learning-paths/dqs0-dmbok-and-cdmp-preparation-data-quality-specialist-learning-plan>

Master & Reference Data Management

<https://www.dataversity.net/>

Metadata Management

<https://www.dataversity.net/>

 CDMP Online learning program

 Approved by DAMA-I

 Part of DATAVERSITY training center

 Based on DMBoK CDMP syllabus

20% discount code for DAMA UK "dmadvisors"

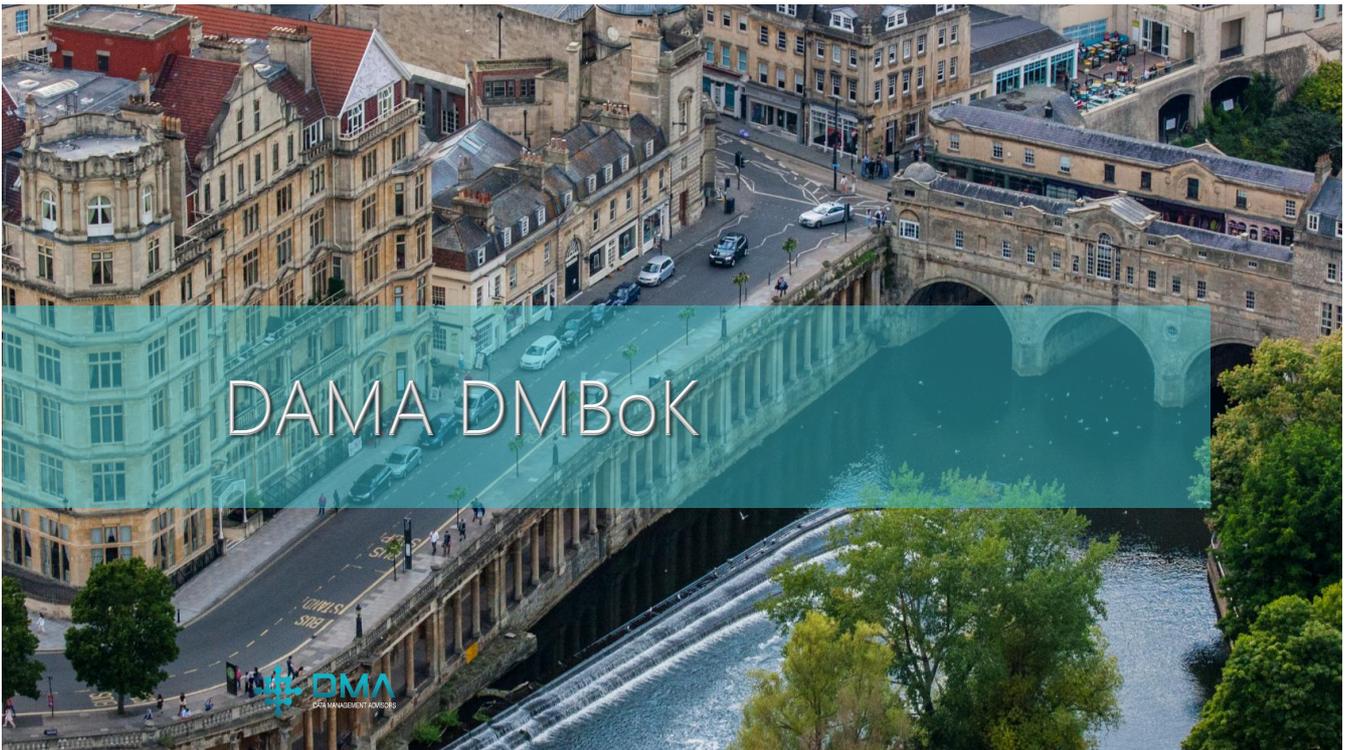


Renewal

To maintain the validity of your Certified Data Management Professional® (CDMP®) certification, DAMA International® requires all certification holders to participate in an annual renewal process. Initially, the CDMP® certification is valid for three years. To keep your certification active beyond this period, you must:

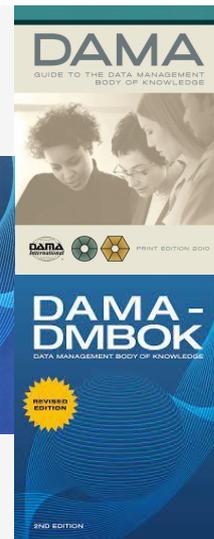
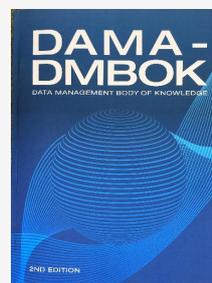
- 1. Annual Maintenance Fee:** Pay the CDMP® maintenance fee each year.
- 2. Data Management Activity Attestation:** Submit the Data Management Activity Attestation Form annually, confirming your ongoing involvement in data management activities.
- 3. Code of Ethics:** Sign and adhere to the current version of the DAMA International® Code of Ethics, reflecting your commitment to professionalism and ethical standards in the field.





What Is the DAMA-DMBOK Guide?

- The DAMA Guide to the Data Management Body of Knowledge (DAMA-DMBOK Guide)
- A book published by DAMA-I, 406 pages(V1) 624 pages (v2) (also on CD & PDF)
- Available from TechnicsPublications.com or Amazon.com
- Written and edited by DAMA members
- An integrated primer: “definitive introduction”
- Modeled after other BOK documents:
 - PMBOK (Project Management Body of Knowledge)
 - SWEBOK (Software Engineering Body of Knowledge)
 - BABOK (Business Analysis Body of Knowledge)
 - CITBOK (Canadian IT Body of Knowledge)



DAMA-DMBOK Guide Goals



Data Management Disciplines

(DMBoK 2)

- Data Management covers considerably more than Documents and Records management
- Data Governance is at the centre
- Data Management is a **professional** discipline which requires training & awareness



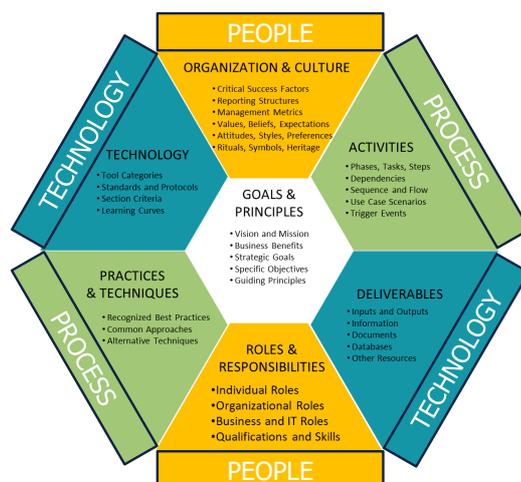
Business Environmental Factors

- Data Management does not exist in a vacuum
- Many environmental factors cross all of the disciplines
- Goals and principles at the centre:
 - Provide guidance how to execute activities and effectively use the tools for successful data management

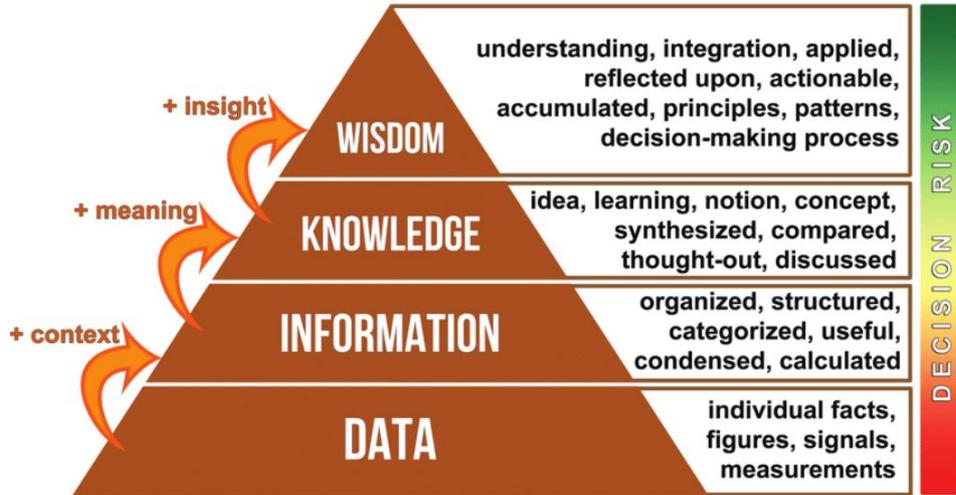


Business Environmental Factors

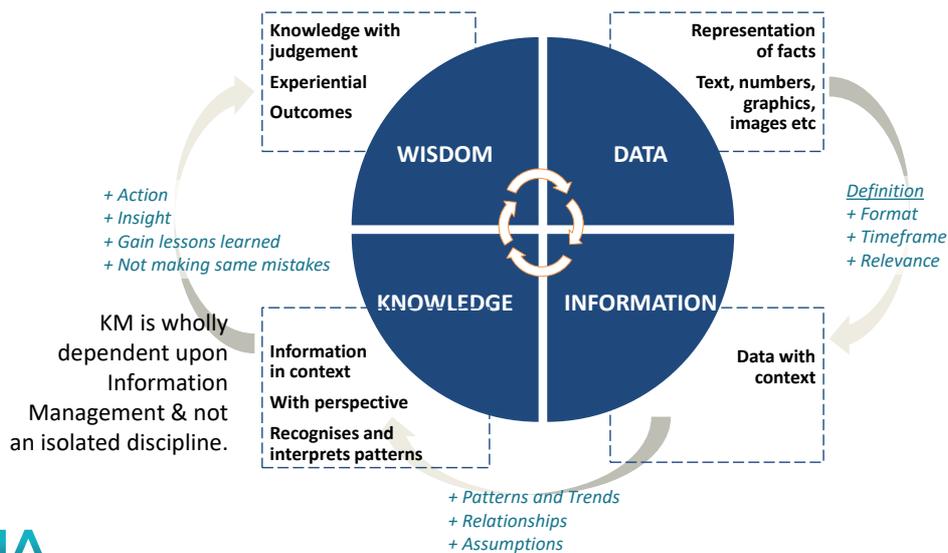
It's NOT just about technology



The Information Value Chain



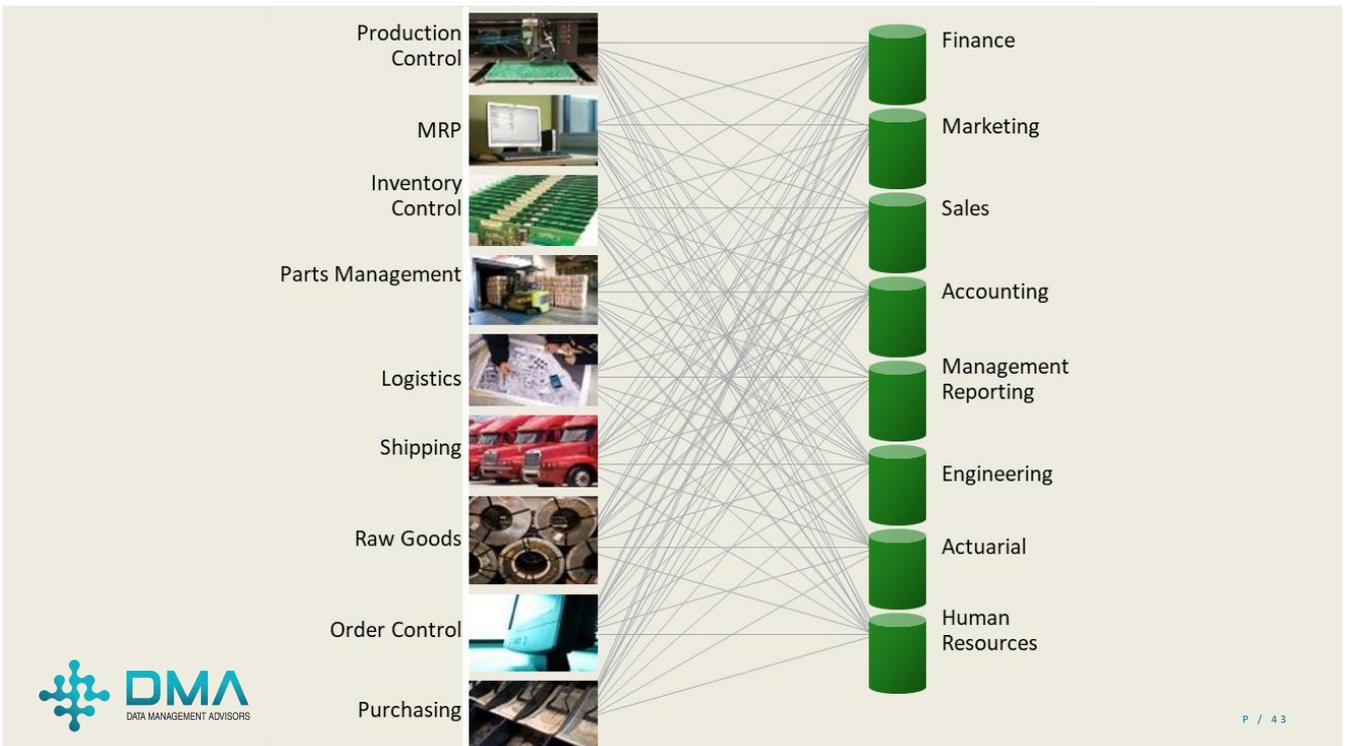
Data, Information and Knowledge *(Information Value Chain)*





Why is Data Management critical?

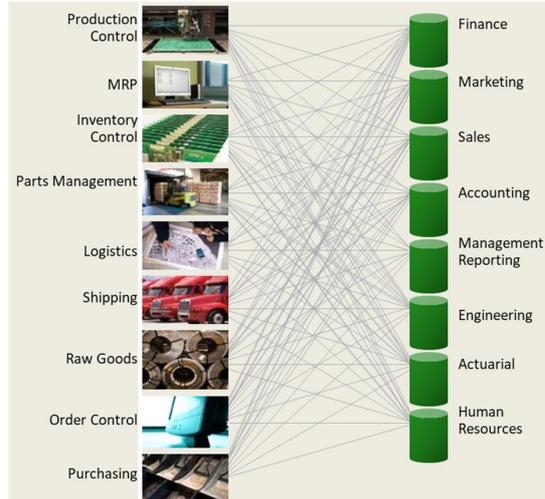
- Higher volumes of data generated by organisations (raw data, devices, CRM, ECM, IOT)
- Proliferation of application-centric systems
- New product development
- To make the management of information front and centre and part of the culture
- Greater demand for reliable information: Gain deep insights through analytics
- Trust in Information: "What do you mean by?"
- Tighter regulatory compliance
- Competitive advantage: Improved decision making
- Business change is no longer optional – it's inevitable: Agility AND ability to respond to change
- Big Data + AI & ML explosion (and hype)



What is the problem?

Legacy Applications Centric focus = **CHAOS**

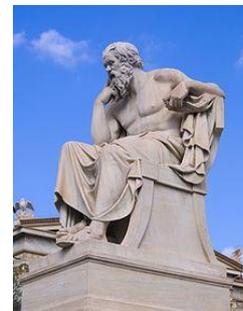
The root cause of the muddled state of Information Architecture in most large organisations today is the prevailing **application-centric** mindset that gives applications priority over data, where data is treated as a “by-product”



Until we recognize and take action on the core problem, the situation will continue to deteriorate.

Clarity & Definition is vital

“The beginning of wisdom is the definition of terms.”



Socrates: 470 – 399 BC



The Problem

+ vendor promises
(with *business willingness to believe*)
in “next new thing”

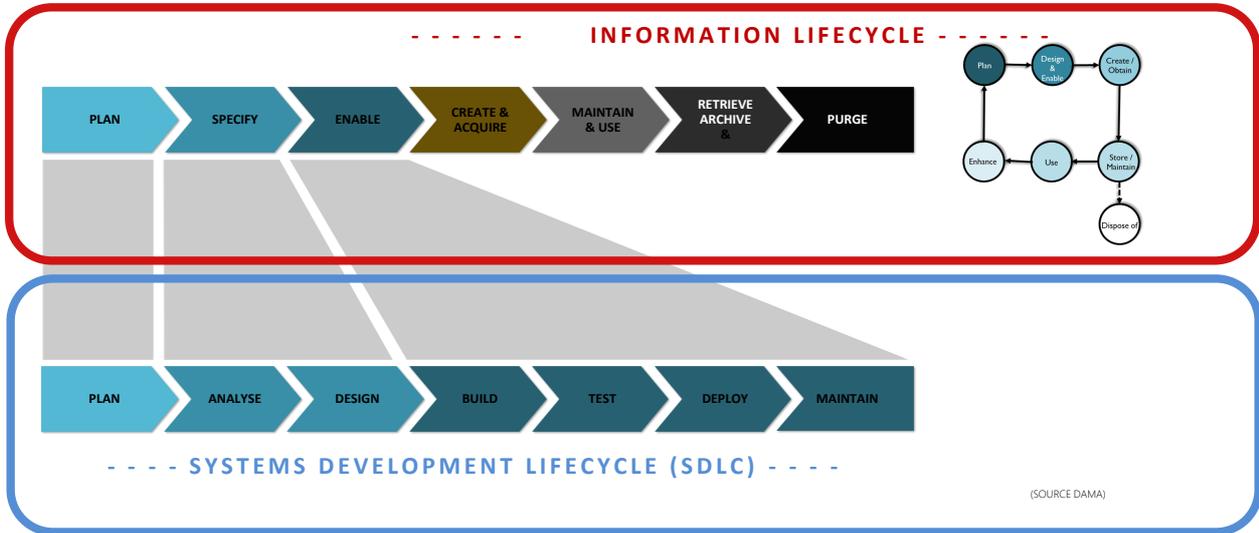


**There's NO
SILVER BULLET**



**They ALL depend on reliable data ... and that's the bit
that's sadly lacking**

Information Lifecycle & SDLC



Data is precious

“Data is a precious thing and will last longer than the systems themselves.”



Sir Tim Berners-Lee:
Inventor of the World Wide Web.

Additional / Optional Study Materials

- DAMA-DMBOK (Technics Publishing) (Required)
- Systems Analysis and Design (Kendall and Kendall - Prentice Hall)
- Master Data Management (Loshin - Morgan Kauffman)
- IT Project Management (Schwalbe - Course Technology)



Ref	Question	A	B	C	D	E
GEN2	Information is	Data in context	A management discipline	Always stored in a computer system	A byproduct of IT systems	An abstract concept
GEN3	Information needs to be managed because	It is an asset of the organization	The volumes are large	It contains financial facts	It is stored in Database systems	Processes use it
GEN4	Data differs with regards to other assets because	It uses automation	It can be used yet still retain value	It has value	It is big	It is regulated
GEN5	The Information Lifecycle	Has the same stages as the Systems Delivery Lifecycle	Is used primarily for Data archiving	Is only important in regulated industries	Exists beyond the Systems Delivery Lifecycle	Is not relevant in an Agile environment
GEN6	The DAMA Wheel contains	Knowledge areas	data management processes	data strategy initiatives	maturity model dimensions	data management deliverables
GEN9	Which is a valid DMBOK Environmental component of data management?	Motivation	Hardware Management	Practices & Techniques	Project Management	Database Management.



Data Quality Management

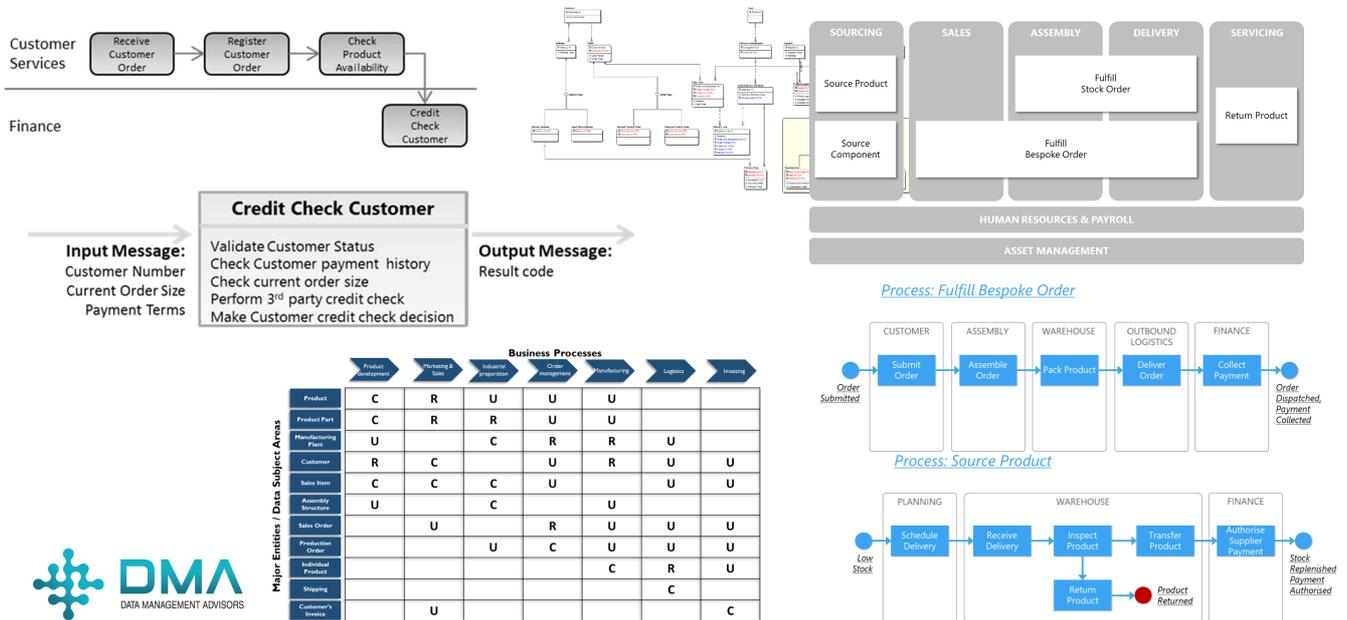
Data Management Process	2%
Big Data	2%
Data Architecture	6%
Document & Content Management	4-5%
Data Ethics	2%
Data Governance	11%
Data Integration & Interoperability	6-7%
Master & Reference Data Management	11%
Data Modelling & Design	11%
Data Quality	11%
Data Security	6%
Data Storage & Operations	4-5%
Data Warehousing & Business Intelligence	11%
Metadata Management	11%



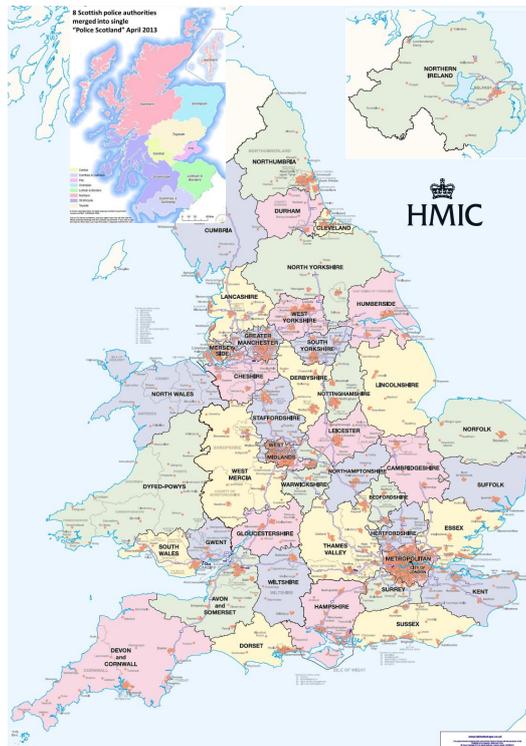
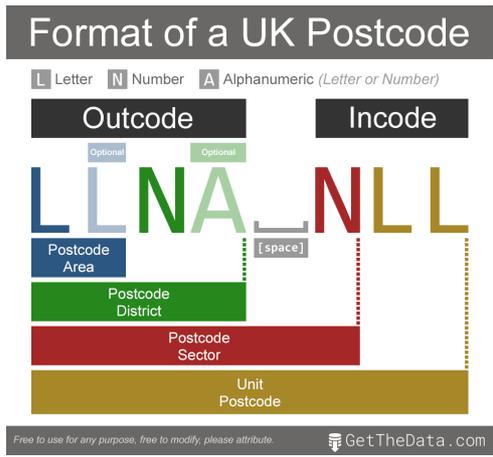
Data Quality Management (DMBoK 2 revised)



Do we know *everything* that the data might be used for?



Data Quality Metrics
 Beware, this is NOT an excuse for omitting requirements analysis !



1 February 2011 Last updated at 14:44

Preston street branded most crime-ridden



A look round 'crime-ridden' Glovers Court, in Preston

A street in Preston has been branded the most crime-ridden in England and Wales by a government website.

Glovers Court - a quiet street in Preston city centre - has been the location for 150 offences in December, according to the online map.

But police say the crimes were actually committed across the whole of Preston city centre.

Ch Supt James Lee said only three crimes had been reported at Glovers Court.

The figures are part of a new online crime map, set up by the Home Office, which cost £300,000 to develop.

Visitors to the website are able to find out which crimes have taken place on or near their street within the past month and which officers are responsible for their area.

The street with the most recorded crimes is Glovers Court - but Ch Supt Lee said: "The figures don't do it justice; it is actually a safe place to be.

I don't accept these figures. The postcode relates to the whole of the city

Related Stories

Millions clog crime map website

Preston street "most dangerous"

Street-level crime maps launched

1 February 2011 Last updated at 16:44

Millions jam street-level crime map website

COMMENTS

A new crime-mapping website for England and Wales is experiencing a "temporary problem" as millions of people log on every hour, the Home Office has said.

Hundreds have contacted the BBC website to report problems accessing the site as officials worked to fix the glitch.

The Home Office said www.police.uk was receiving up to five million hits an hour, or some 75,000 a minute.



The maps record reports of different categories of crime and anti-social behaviour

The site allows you to see the offences reported in your local street by entering a street name or postcode.

Home Secretary Theresa May said the maps would give real facts on crime and anti-social behaviour and make police more accountable.

'Complete farce'

In a message on the microblogging website Twitter, the Home Office said: "Hugely popular streetlevel crime maps getting 75,000 hits per minute so you might experience delays. Keep trying."

But some users have reported seeing only a blank page when visiting the website, or a message saying "no police area is associated with this address" when entering a postcode or street.

Others have complained of errors in the actual information shown, with some "quiet streets" next to commercial centres, bars and clubs being tarred with their crimes.

A spokesman for charity Victim Support said it was important that victims of crime had consented as to whether information about their incident was released.

Related Stories

Street branded most crime-ridden

Street-level crime maps launched

Patrolling a city's flashpoints



Seeing the dots relating to 1,672 incidents a mile from my north London front door was quite sobering"

Mark Easton

P / 59

How Good Does Data Quality Need To Be?

In February 2011, the UK government launched a crime-mapping website for England and Wales (www.police.uk).

Unfortunately, for a number of reasons, the postcode allocated to a specific police incident didn't always correspond to the precise location of the crime.

The net result was that poor accuracy in the recording of geographical information led many quiet residential streets to be incorrectly identified as crime hotspots.

Answer: It depends...



In the context of creating aggregated statistics to assess relative crime rates between counties, the data quality is perfectly acceptable.

✓ Data fit for purpose

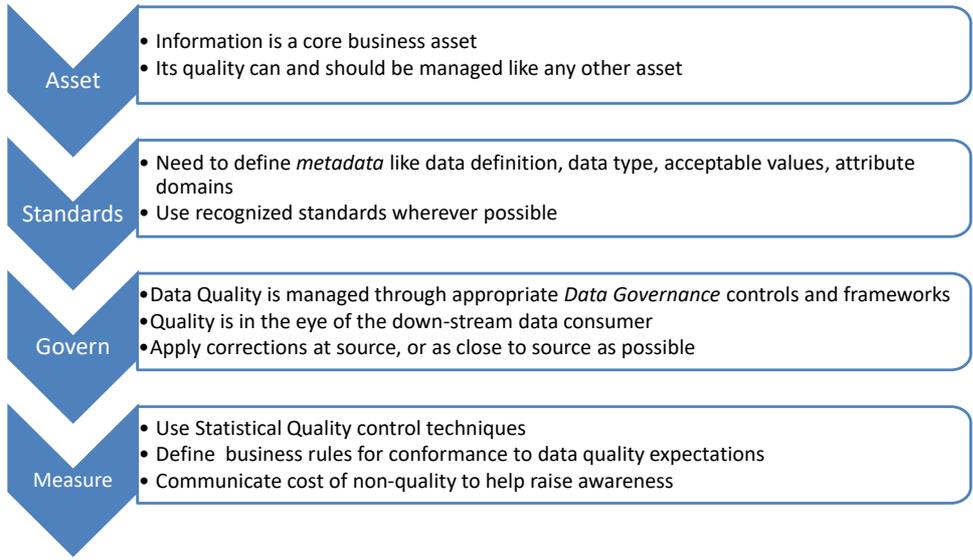


However, if the same data is used by an insurance company, there is an issue for the homeowners who receive inflated home insurance premiums.

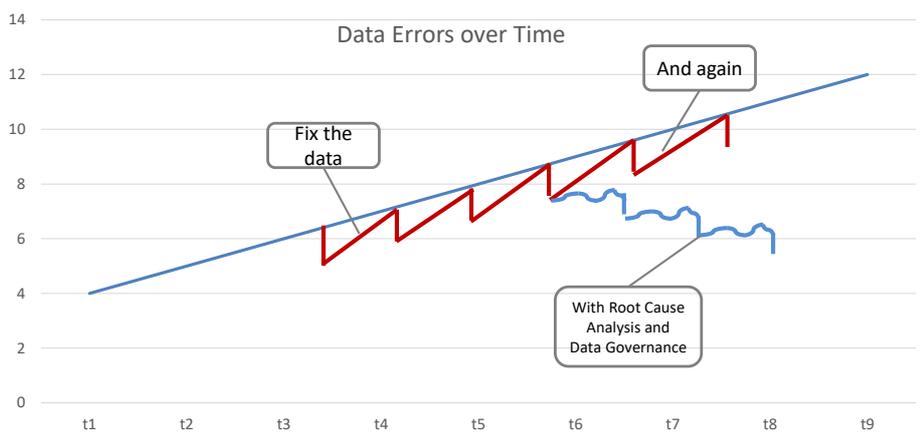
✗ Data not fit for purpose

Data quality can only be considered within the context of the *intended use of the data*
Data needs to be "fit for purpose"
Data quality needs to be assessed on that basis

Key Points



Don't just fix the data



Data Quality Principles (DMBoK 2)



Criticality: A Data Quality improvement program should focus on the data most critical to the enterprise and its customers. Priorities for improvement should be based on the criticality of the data and on the level of risk if data is not correct.



Standards-driven: All stakeholders in the data lifecycle have Data Quality requirements. To the degree possible, these requirements should be defined in the form of measurable standards and expectations against which the quality of data can be measured.



Objective measurement and transparency: Data Quality levels need to be measured objectively and consistently. Measurements and measurement methodology should be shared with stakeholders since they are the arbiters of quality.



Prevention: The focus of Data Quality improvement program should be on preventing data errors and conditions that reduce the quality of data; it should not be focused on simply correcting records.



Root cause remediation: Improving the quality of data goes beyond correcting errors. Problems with the quality of data should be understood and addressed at their root causes rather than just their symptoms. Because these causes are often related to process or system design, improving Data Quality often requires changes to processes and the systems that support them.



Embedded in business processes: Business process owners are responsible for the quality of data produced through their processes. They must enforce Data Quality standards in their processes.



Systematically enforced: System owners must systematically enforce Data Quality requirements.



Connected to service levels: Data Quality reporting and issues management should be implemented and incorporated into Service Level Agreements (SLA).



Critical Data

Critical Data Elements are identified by linking them to the Data Quality business drivers; customer experience, effectiveness, and efficiency.

Critical data needs to be managed to a defined Data Quality to meet these objectives.

- Without a particular data element there will be significant impairment to the success of the organization. This impairment size drives a business case for the level of Data Quality Management of a particular data element.
- This value prioritization approach creates a sliding scale of value.

At a defined value level, data elements will become “critical data elements”.

- EG if the data in the customer email address field is incomplete, we will not be able to send product information to our customers via email and we will lose potential sales.
- We know that for every email we send out, we earn \$1.00 in revenue.
- This simple example links a clear value driver to Data Quality improvement.

Critical data is often used in:

- Regulatory, financial, or management reporting
- Business operational needs
- Measuring product quality and customer satisfaction
- Business strategy, especially efforts at competitive differentiation.

Master and Reference Data is usually critical.

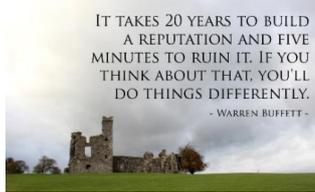
- Data sets or individual data elements can be assessed for criticality based on:
 - the processes that consume them,
 - the nature of the reports they appear in,
 - the financial, regulatory, or reputational risk to the organization if something were to go wrong with the data.



The Impact of Poor Data



AGGRAVATION



REPUTATION



TELL THEIR FRIENDS



DEATH



RISK



LOSS



<https://www.manchestereveningnews.co.uk/news/greater-manchester-news/cracked-gas-main-likely-caused-21353607>

<https://www.bbc.co.uk/news/av/uk-england-12335044>

<https://www.theguardian.com/business/2019/may/22/record-44m-fine-for-cadent-after-residents-had-no-gas-for-months>

<https://www.bbc.co.uk/news/uk-england-lancashire-35048851>

<http://news.bbc.co.uk/1/hi/scotland/4184962.stm>

<https://ico.org.uk/about-the-ico/media-centre/news-and-blogs/2020/10/ico-fines-british-airways-20m-for-data-breach-affecting-more-than-400-000-customers/>

[British Airways fined £20m over data breach - BBC News](#)

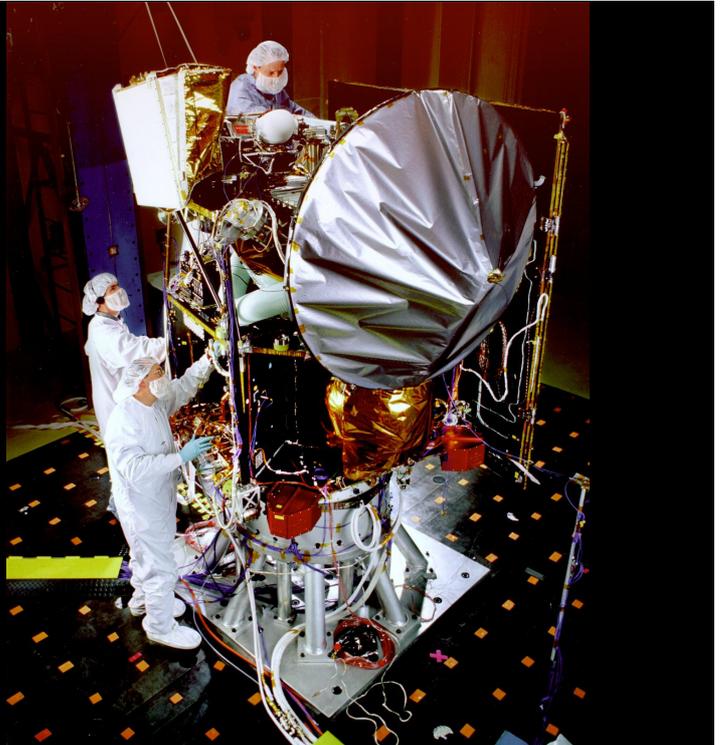
[Google AI chatbot Bard sends shares plummeting after it gives wrong answer | Google | The Guardian](#)

Data being used for the wrong thing: <https://www.abc.net.au/news/2023-02-20/robodebt-scheme-government-royal-commission-into-fraud-income/101998782>

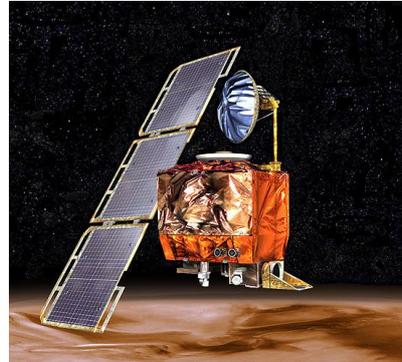
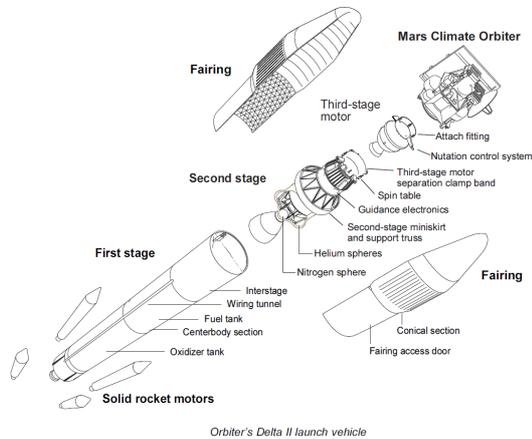




Launch date: 11 December 1998



Mars Climate Orbiter

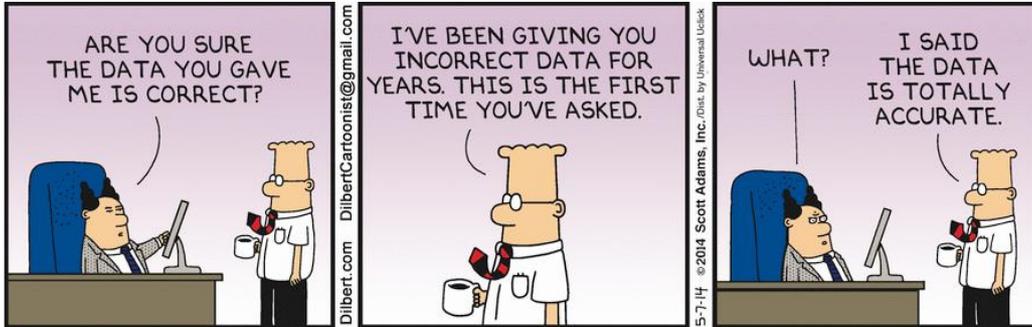


What Do You Mean By That?

Decay date: 23 September 1999, "Unintentionally deorbited"

Ground software supplied by Lockheed Martin produced results in **pound seconds** ("American"), contrary to its Software Interface Specification (SIS), while a second system, supplied by NASA, used the results expecting them to be in **newton-seconds**, as per the SIS.

The discrepancy between calculated and measured position, resulting in the discrepancy between desired and actual orbit insertion altitude, had been noticed earlier by at least two navigators, *whose concerns were dismissed*.

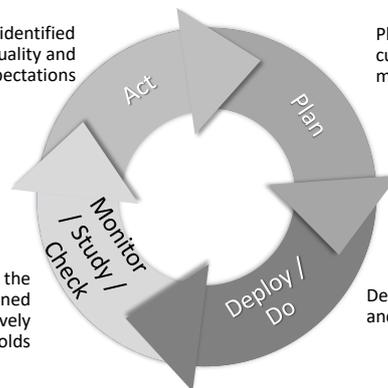


Data Quality Management Cycle

The Data Management Body of Knowledge identifies 4 key activities necessary for operationalising DQM:
Plan, Deploy, Monitor, Act

Acting to resolve any identified issues to improve data quality and better meet business expectations

Planning for the assessment of the current state and identification of key metrics for measuring data quality



Monitoring and measuring the levels in relation to the defined business expectations & actively checking data against thresholds

Deploying processes for measuring and improving the quality of data

The Data Quality Management Approach

Shewart Cycle	Deming Cycle	DQM Cycle	Summary of Tasks
Plan	Plan	Plan	What data issues are critical to achievement of business objectives? What are requirements for data quality? What are dimensions of DQ and business rules
Do	Do	Deploy	Conduct data profiling, inspect and monitor for data quality issues; Identify root causes of data quality issues; Apply process or data remediation at, or near, source; Work to eliminate “special causes” of defect
Check	Study	Monitor	Implement active monitoring of DQ. Apply statistical control. Investigate and understand defects that arise, make plans for remediation
Act	Act	Act	Take actions to remedy emerging data quality issues as they are identified.

Juran Trilogy

(Joseph Juran 1908-2008)

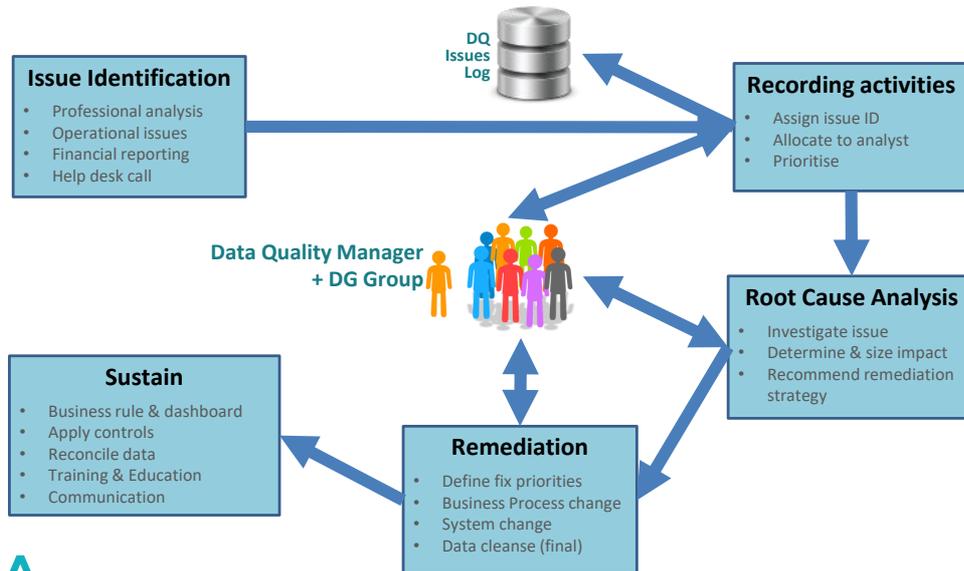


Quality Planning	Quality Control	Quality Improvement
Establish Quality Goals Identify who the customers are Determine the needs of the customers Develop product features that respond to customers needs Develop processes able to produce the product features Establish process control; transfer the plan to the operating forces	Evaluate Actual Performance Compare Actual Performance with Quality goals Act on the difference	Prove the needs Establish the infrastructure Identify the improvement projects Establish the project teams Provide the teams with resources, training and motivation to diagnose the causes and stimulate remedies Establish controls to hold the gains

Underlying concept: “managing for quality consists of three universal processes”

- **Quality Planning (Quality by Design)**
- **Quality Control (Process Control & Regulatory)**
- **Quality Improvement (Lean Six Sigma)**

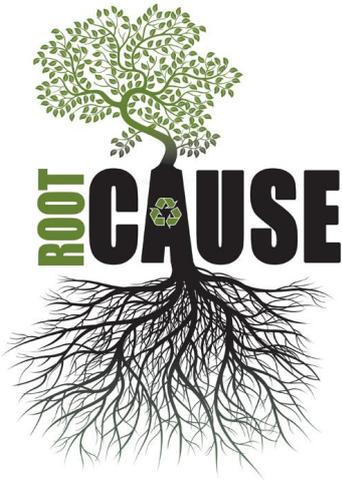
A simple Data Quality improvement framework



DQ Governance

The steps followed in managing data issues include:





- Step 1:** Define the Problem
- Step 2:** Collect Data
- Step 3:** Identify Possible Causal Factors
- Step 4:** Identify the Root Cause(s)
- Step 5:** Recommend and Implement Solutions

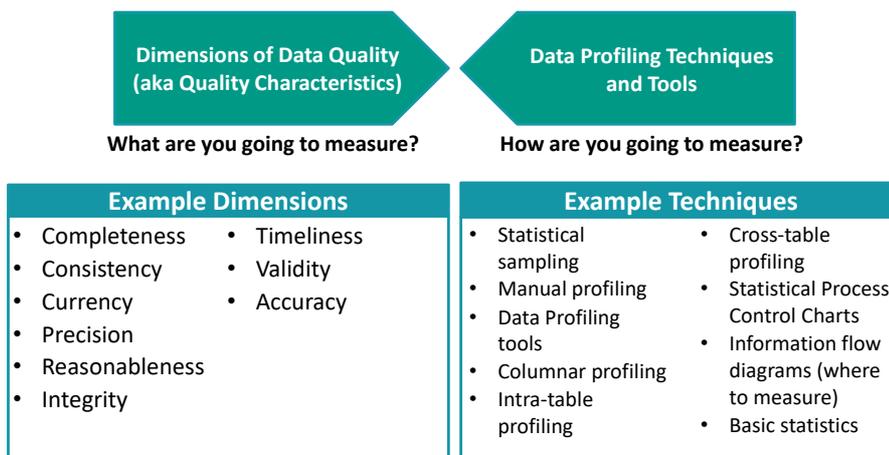


MAN: HUMAN CAUSES – People did something wrong, or did not do something that was needed. Human causes typically lead to physical causes (for example, no one filled the brake fluid, which led to the brakes failing).

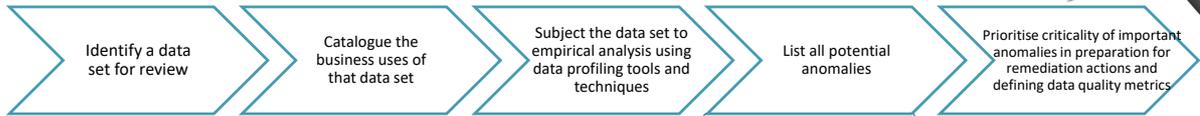
METHOD: ORGANIZATIONAL CAUSES – A system, process, or policy that people use to make decisions or do their work is faulty (for example, no one person was responsible for vehicle maintenance, and everyone assumed someone else had filled the brake fluid).

MACHINE: PHYSICAL CAUSES – Tangible, material items failed in some way (for example, a car's brakes stopped working).

Measuring Data Quality

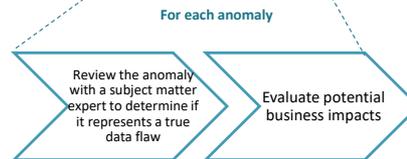


Data Quality Profiling – Assess Data Quality



Typical outputs from Data Quality Profiling

- % of the records populated.
- The number of data values populating each data attribute.
- Frequently occurring values.
- Potential outliers.
- Relationships between columns within the same table.
- Relationships across tables



Why Data Quality Profiling?

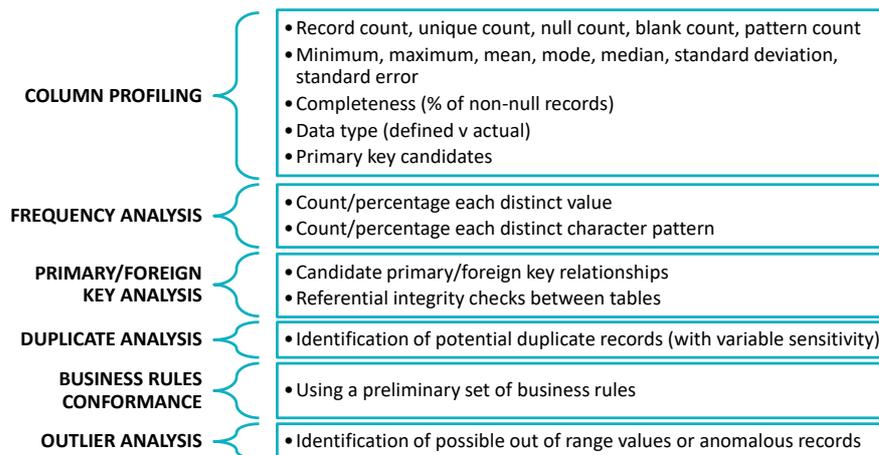
- Reviewing and refining business policies provides a “top down” view of Data Quality requirements, but a “bottom up” view is crucial to identify existing issues within the data
- This is achieved through an activity known as Data Quality Profiling
- To conduct Data Quality Profiling as efficiently and repeatedly as possible, a specialist Data Quality tool is normally employed
- The result is an invaluable insight into the real operational data, revealing hidden characteristics, patterns and anomalies

Activities:

1. Define High Quality Data (P)
2. Define a Data Quality Strategy (P)
3. Define Scope of Initial Assessment (P)
 1. Identify Critical Data
 2. Identify Existing Rules and Patterns
4. Perform Initial Data Quality Assessment (P)
 1. Identify and prioritize issues
 2. Perform root cause analysis of issues
5. Identify & Prioritize Improvements
 1. Prioritize Actions based on Business Impact
 2. Develop Preventative and Corrective Actions
 3. Confirm Planned Actions
6. Develop and Deploy Data Quality Operations (D)
 1. Develop Data Quality Operational Procedures
 2. Correct Data Quality Defects
 3. Measure and Monitor Data Quality
 4. Report on Data Quality levels and findings



Typical Outputs of Data Quality Profiling



Data Profiling



experian. Aperture Data Studio Space: Demo

Customer List Exploring

Filter Sample Columns Sort Group Transform Infer Profile Download CSV Save as View Discard

	Name	Uniqueness	Unique Count	Completeness	Row Count	Has Nulls	Overall Datatype	Dominant Datatype	Format Count	Most Common Format	Most Common Value	Longest Length	Shortest Length
1	Customer Id	92.5%	185	100%	200		Alphanumeric	Alphanumeric	1	A99-A999999	A89-E137529	11	11
2	Discount Code	2%	4	100%	200		Alphanumeric	Alphanumeric	1	A	N	1	1
3	Forename	66.5%	133	100%	200		Alphanumeric	Alphanumeric	15	AAAA	Accounts	11	1
4	Surname	83%	166	95%	200		Alphanumeric	Alphanumeric	17	AAAAA		35	3
5	Email	94%	188	100%	200	Yes	Alphanumeric	Alphanumeric	182	AAA_AAAAAAAAAA@	Lee_Reichert@Wpico.	58	19
6	Telephone	94%	188	100%	200		Alphanumeric	Alphanumeric	11	(999)999-9999	(973) 467-8780	16	9
7	Company Name	94%	188	100%	200		Alphanumeric	Alphanumeric	165	AAAAASAAAAAAA	WPI Communications.	86	6
8	First Order Date	94%	188	100%	200		Alphanumeric	Date	3	99-AAA-9999	13-Sep-1972	11	9
9	Sales Last Year	66%	132	100%	200		Number	Number	5	99.999	10,000	6	1
10	Addressline	95.5%	191	100%	200		Alphanumeric	Alphanumeric	171	AASAAAS99999	2900 Telestar Cr.	40	8
11	City	68%	136	100%	200		Alphanumeric	Alphanumeric	39	AAAAAA	New York City	22	3
12	Country	0.5%	1	100%	200		Alphanumeric	Alphanumeric	1	AAAAASAAAAAAA	United States	13	13
13	Post Code	93.5%	187	100%	200		Alphanumeric	Alphanumeric	6	99999-9999	07081-1426	10	4
14	Good Orders	87.5%	175	100%	200		Number	Number	3	999	187	3	1
15	Return Orders	41%	82	100%	200		Number	Number	3	99	7	3	1
16	Total Orders	87.5%	175	100%	200		Number	Number	4	999	206	4	1

16 rows 65 columns



Column Level Stats 1

Summary Info

experian. Aperture Data Studio Space: Demo

Customer List Exploring

Filter Sample Columns Sort Group Transform Infer Profile Download CSV Save as View Discard

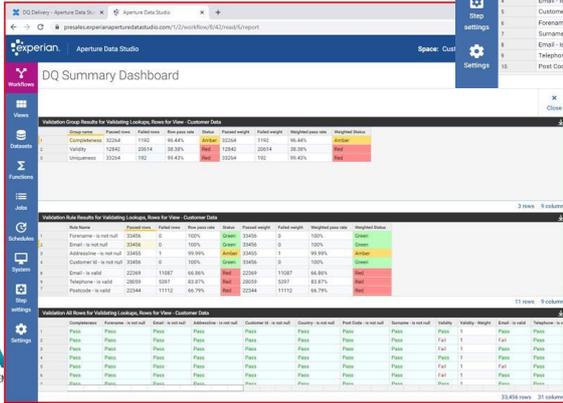
	Value (Surname)	Row count	Distribution	Length	Format	Simple format	English phonetic	Soundex	Refined soundex	Datatype
1		10	5	0						NULL
2	Green	4	2	5	AAAAA	A	KRN	G650	G4908	ALPHANUMERIC
3	Reichert	4	2	8	AAAAA	A	RXRT	R263	R903096	ALPHANUMERIC
4	Kalmakov	3	1.5	8	AAAAA	A	KLMKF	K452	K30780302	ALPHANUMERIC
5	Roberts	3	1.5	7	AAAAA	A	RPRTS	R163	R9010963	ALPHANUMERIC
6	Taylor	3	1.5	6	AAAAA	A	TLR	T460	T60709	ALPHANUMERIC
7	Williams	3	1.5	8	AAAAA	A	ALMS	W452	W07083	ALPHANUMERIC
8	Ames	2	1	4	AAAA	A	AMS	A520	A0803	ALPHANUMERIC
9	Barker	2	1	6	AAAAA	A	PRKR	B626	B109309	ALPHANUMERIC
10	Brown	2	1	5	AAAAA	A	PRN	B650	B1908	ALPHANUMERIC
11	Grantham	2	1	8	AAAAA	A	KRNTM	G653	G4908608	ALPHANUMERIC
12	Lloyd	2	1	5	AAAAA	A	LT	L300	L706	ALPHANUMERIC
13	Morgan	2	1	6	AAAAA	A	MRKN	M625	M809408	ALPHANUMERIC
14	Parry	2	1	5	AAAAA	A	PR	P600	P1090	ALPHANUMERIC
15	Thompson	2	1	8	AAAAA	A	TMPSN	TS12	T6081308	ALPHANUMERIC
16	Turner	2	1	6	AAAAA	A	TRNR	T656	T609809	ALPHANUMERIC
17	Whitmore	2	1	8	AAAAA	A	ATMR	W356	W068090	ALPHANUMERIC

166 rows 10 columns

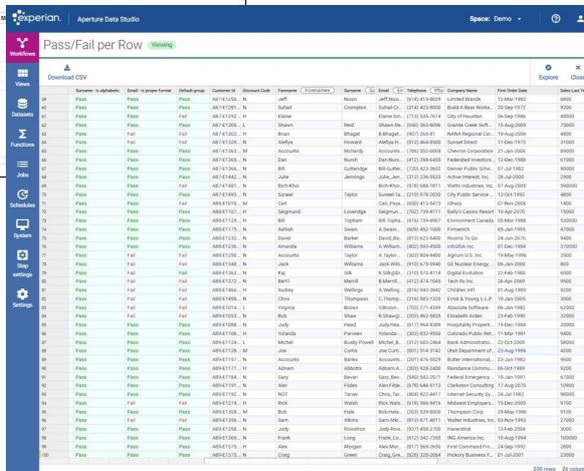
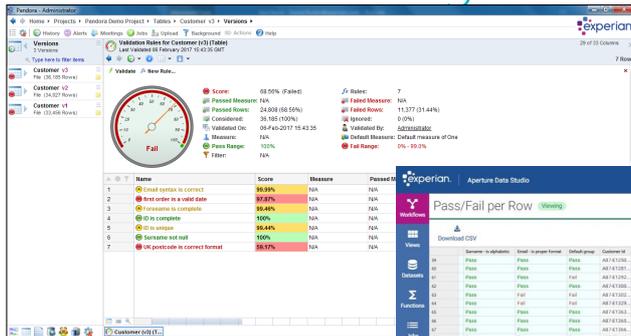


Validation rules

- Every evaluated column can be scored
- Every failing and passing row can be explicitly identified
- Every failure reason is clearly identifiable



Rules-based Monitoring



Data Quality Tools and Techniques

Technique	Description	Use
Data Profiling	An objective review of data values against expected standards	Use to measure DQ levels and identify data that is "fit for purpose"
Parsing & Standardisation	Use of patterns and rules to extract data into discrete tokens or to apply standard values based on a defined reference standard	Use to extract data into more granular formats or to identify components in strings; Reduce variation in data set
Transformation	Use of triggered data rules to transform source data to target structures	Use to cleanse/rework data from source to target. Builds on data standardisation rules
Identity resolution & matching	Clustering of records based on "similarity" to identify duplicates or cluster households	Use to dedupe databases or identify important clusters
Enhancement (aka enrichment)	Adds value to data by linking additional information to a record or correcting errors in a record when compared against a definitive source. Can be used to add detail to the record	Use to increase the value of the data, e.g. by appending a market segment to a customer record. Other examples include appending metadata to records to support data lineage tracking or other contextual data.

Dimensions of Data Quality



Search on Google & you'll see lots of "Dimensions"

- Experian / DAMA UK & most commentators = 6
- EDM Council / Acuate DQ = 7
- David Loshin / others = 8
- DMBOK (V2) = 11
- DMBOK (V2Rev) = 9

Data Quality Dimensions DMBok2 Revised Edition

Quality Attribute	Description
Accuracy	Equal to the “real life” value
Completeness	All mandatory values present, all optional, but implied values present
Consistency	A value in one data set aligns to another data set
Currency	How “fresh” is the data compared to the real-world
Precision	The level of detail in the data element – covered in Validity & Accuracy
Privacy	Does the data need restricted access or monitoring? – covered in Metadata classifications
Reasonableness	Is the current behaviour in line with previous behaviour (or average)?
Referential Integrity	All child records must have a parent.
Timeliness	The difference between when the data is available and needed
Uniqueness	Are there duplicates within the data (that do not reflect reality!)
Validity	Is the data within the allowable bounds of the “domain”.

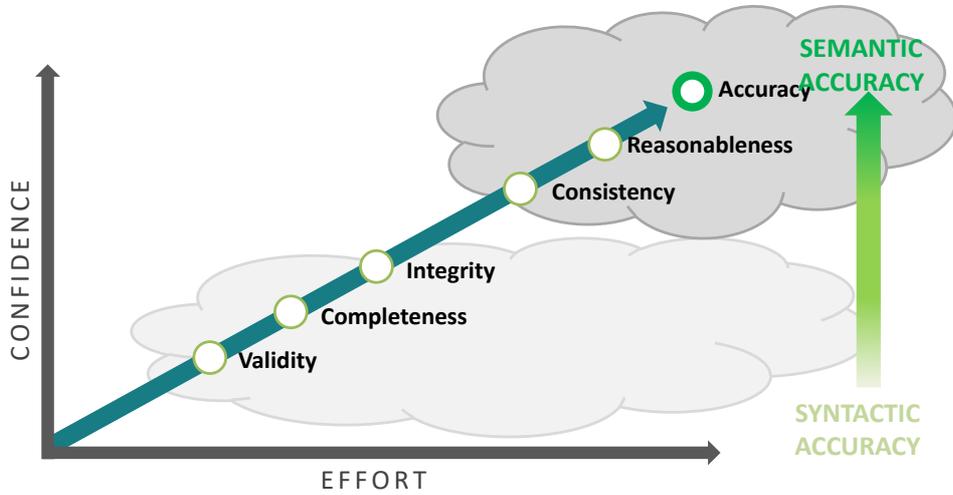


DIM	DESCRIPTION	EXAMPLES
Validity	Validity refers to whether data values are consistent with a defined domain of values.	<p>A domain of values may consist of:</p> <ul style="list-style-type: none"> • A data type (e.g., text with or without special characters, a number, or a date) • A data range (e.g., a numeric range, date range, or a set of valid text values in a reference table) • A format (e.g., a phone number must contain country, area codes, and the correct number of digits or a currency must have only two decimal places) • The precision of expected values. For example, a date-stamp must be recorded to a millisecond or a number recorded to no more than two decimal places. • The time frame (certain values may also only be valid for a specific length of time, for example, data that is generated from RFID (radio frequency ID) or some scientific data sets).
Completeness	<p>Completeness refers to whether all the required data is present. For example, are all the mandatory components of an address populated, including the street number, street name, city, and country?</p> <p>For example, if the customer account is marked as a debt risk, there must be a debt manager appointed.</p> <p>If the customer account is not marked as a debt risk, there must not be a debt manager appointed.</p>	<p>Completeness can be measured at the column, record, or data set level.</p> <ul style="list-style-type: none"> • Are columns/data elements populated to the level expected? (Some columns are mandatory. Optional columns are populated only under specific conditions.) • Are records populated correctly? (Records with different statuses may have different expectations for completeness.) • Does the data set contain all the records expected?
Consistency	<p>Consistency is ensuring that data values are coded using the same approach, assessment and valuation criteria.</p> <p>Consistency is between two different data items, potentially within a data set and between data sets, and across time.</p> <p>Consistency may be defined between one set of data element values and another data element set within the same record (record-level consistency), between one set of data element values and another data element set in different records (cross- record consistency), or between one set of data element values and the same data element set within the same record at different points in time (temporal consistency).</p> <p>Consistency can also be used to refer to consistency of format.</p> <p>Take care not to confuse consistency with accuracy or correctness.</p>	<p>Consistency can be measured using the following ways:</p> <ul style="list-style-type: none"> • Cross record consistency is between data values in the same column – are all the customer addresses recording the head-office details or do some record service delivery centres? • Between associated data sets, where data sets are expected to have values that are linked. For example, when moving data from data capture systems to data warehousing systems, a consistency quality measure will maintain confidence in the values in the data warehouse. • Across time can be applied to any of the previous cases. For example, has the grading of students been the same over time? If a subject is completed previously, would the same learning performance yield the same grade today?

DIM	DESCRIPTION	EXAMPLES
Integrity	Integrity refers to the lack of incoherent values and broken relationships in data.	<p>Integrity can be measured using the following ways:</p> <ul style="list-style-type: none"> Coherence – where one data value implies a limited range in another data value, and they are matching. For example, a country has a defined set of states or provinces. If a customer address is in a certain country, they must use a state or province name for the associated set. For example, Alberta is a province of Canada, Pisa is a province of Italy, and Tasmania is a state of Australia. Parent-child (referential integrity) – where every child must have a parent data value. For example, all the customer address country names must exist on the valid list of countries where the organization is authorized to sell products to. <p>Data sets without integrity are seen as corrupted or have data loss. Data sets without referential integrity have ‘orphans’ – invalid reference keys, or ‘duplicates’ – identical rows which may negatively affect aggregation functions.</p>
Timeliness	Data Timeliness refers to the time that data needs to become accessible to a user after its capture or update. Timeliness is the expectation and actual delay before the data becomes available.	<p>For example, An electricity utility needs to maintain its power grid and must have demand data measured and available to a system operator within a few seconds.</p> <p>Otherwise, electricity generation may get out of alignment with demand, causing an energy excess or shortage. As another example, a statistical agency of government produces the Gross Domestic Product (GDP) report two months after the end of the quarter. To meet this timeliness requirement, all the data collections that make up the GDP will have their own specific timeliness requirements.</p>
Currency	<p>Currency is the date that the data was last updated relative to now and the likelihood that it is still correct. Different data sets will have different currency expectations from relatively static data to highly volatile data.</p> <p>Static data remains current for a long period. Volatile data remains current for a relatively short period.</p> <p>For example, country codes are relatively static, remaining current for a long period. A business’ bank account is relatively volatile as the balance is constantly changing.</p>	<p>Volatile data, like stock prices on financial web pages, will often be shown with an “as-of-time”, so that data consumers understand the risk that the data has changed since it was recorded. During the day, while the markets are open, such data will be updated frequently – it is volatile. Once markets close, the data will remain unchanged but will still be current since the market itself is inactive –it is static.</p> <p>Key terms used in currency include:</p> <ul style="list-style-type: none"> Update time – the time stamp of the last update Volatility measures the rate of change of the date. It may measure across all the data values or measure within a segment. For example, the rate that domestic customers change their address may be different from the rate that corporate customers change their address. Latency measures the time between when the data was created and when it was made available for use. For example, overnight population of a data warehouse may end up with different data latencies. Data can have a latency of 1 day for data entered into the system early on the prior day, but only a few minutes for data generated just before the load.

DIM	DESCRIPTION	EXAMPLES
Reasonableness	<p>Reasonableness asks whether a data pattern meets expectations.</p> <p>For example, whether a distribution of sales in a geographic area makes sense based on what is known about the customers in that area.</p>	<p>For example, based on previous customer logins at 5pm, are today’s customer logins to our systems out of the ordinary?</p> <p>Key terms used in reasonableness include:</p> <ul style="list-style-type: none"> Fixed values – this is where the measure is fixed, not dependent on previous behaviour. The maximum currency value of a transaction may be limited in size by the technical limits of a bank’s systems and other values bigger than this value is “unreasonable”. Benchmark value – it is possible to measure the average, deviation or any other statistical function and ensure the latest set of values conforms to expected boundaries. For example, there have been between 5,000 and 10,000 customer address changes per day, but yesterday there were 2,000,000. On a government data collection, the average size company has 20 employees, but the latest returns had an unusual number with an average of 1,500 employees reported.
Uniqueness	<p>Uniqueness states that no real-world entity exists more than once within the data set.</p> <p>Asserting uniqueness of the entities within a data set implies that each row relates to each unique entity in the real world, and only that specific entity.</p>	<p>Uniqueness can be identified using the following techniques:</p> <ul style="list-style-type: none"> Key Structure – where there are duplicate keys within the data set. For example, in a data set containing customers, many of which share the same customer number. Related data – where other data in the data may point to a duplicate. For example, two customers with different customer numbers, but having the same name, date of birth, and office address. Inquiries or complaints that users cannot access certain data. For example, a customer with multiple services can only access some of their services as the others are recorded under a different customer number (which is a duplicate).
Accuracy	<p>Accuracy refers to the degree that data correctly represents ‘real-life’ entities.</p> <p>For example, is the person’s name in our database actually the person’s name in real life?</p> <p>Does the customer actually use that email address?</p> <p>Accuracy is difficult to measure unless an organization can reproduce data collection or manually confirm accuracy of records.</p> <p>Most measures of accuracy use the other dimensions to imply accuracy.</p>	<p>Common techniques to improve accuracy include:</p> <ul style="list-style-type: none"> Checking consistency with a data source that has been verified as accurate; a government company register (system of record) or a commercial company Information Broker (system of reference). Ongoing calibration of devices with reality. For example, over time, weather gauges may drift from providing accurate measurements. An annual program of field inspections visits 10% of the devices and recalibrates them to a standard to ensure they are delivering accurate measurements. Sampling data values with reality. For example, small number of emails are sent to customers (mainly) to drive sales, but also have the benefit of checking the accuracy of the email address (does it exist, is it used)?

The Path to Accuracy



ISO 8000 & 22745



ISO 8000 defines quality data as *"portable data that meets stated requirements."* The Data Quality standard is related to the ISO's overall work on data portability and preservation.



Data is considered *'portable'* if it can be separated from a software application. Data that can only be used or read using a specific licensed software application is subject to the terms of the software license.



An organization may not be able to use data it created *unless that data can be detached from the software* that was used to create it.



To meet stated requirements requires that these requirements be defined in a clear, unambiguous manner.



ISO 8000 is supported through ISO 22745, a standard for defining and exchanging Master Data.



ISO 22745 defines how data requirement statements should be constructed, provides examples in XML, and defines a format for the exchange of encoded data.



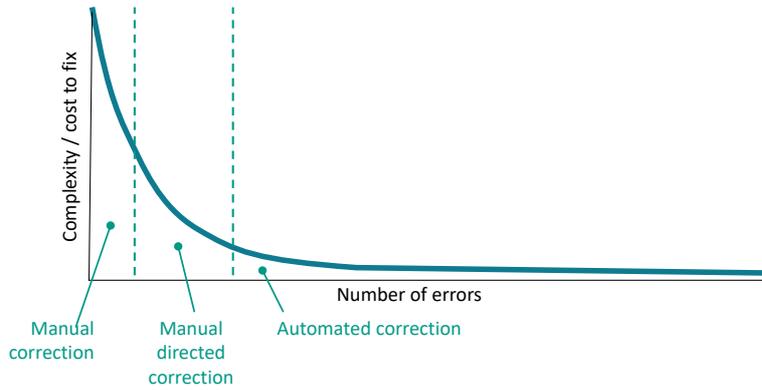
ISO 22745 creates portable data by labelling the data using an ISO 22745 compliant Open Technical Dictionary such as the Electronic Commerce Code Management Association (ECCMA) Open Technical Dictionary (eOTD).



Data Cleansing Demystified

The *quic*k fox jump's over the the lazy dog

Data Correction



Automated correction:

Data transformations and rule-based standardizations, normalizations, reference tables, and corrections. Modified values are committed *without* manual intervention.

Manual directed correction:

Automated tools cleanse and correct data, up to a point. Requires *manual* review before committing the corrections to persistent storage. E.g. scoring mechanism is used to propose a level of confidence in the correction & corrections with scores above a particular level of confidence may be committed without review

Manual correction:

SMEs & Data stewards inspect invalid records and determine the correct values, make the corrections, and commit the updated records. E.g. Pharmaceutical compound & Clinical trial data



Measures at 3 levels of granularity

Granularity	Common Dimensions	Treatment
Data Element	Validity, Integrity, Consistency, Reasonableness	Edit checks in application Data element validation services Specially programmed applications
Data Record	Consistency, Completeness, Currency	Edit checks in application Data record validation services Specially programmed applications
Data set	Completeness, Uniqueness, Reasonableness	Inspection inserted between processing stages



Caution When Defining Data Quality Indicators

DEFINE DATA QUALITY INDICATORS (DQI) WITH THESE CHARACTERISTICS:

Trackability – make sure each DQI is monitored over time to track progress (NB DMBOK uses the term TRENDING)

Acceptability – make sure it’s possible to define what “good” looks like for each DQI – baseline to “goal” = 100%

Relevance – make sure each DQI measures something of importance to the business

Measurability – make sure each DQI can be measured and quantified

Accountability/Stewardship – make sure each DQI links to the Data Governance structure

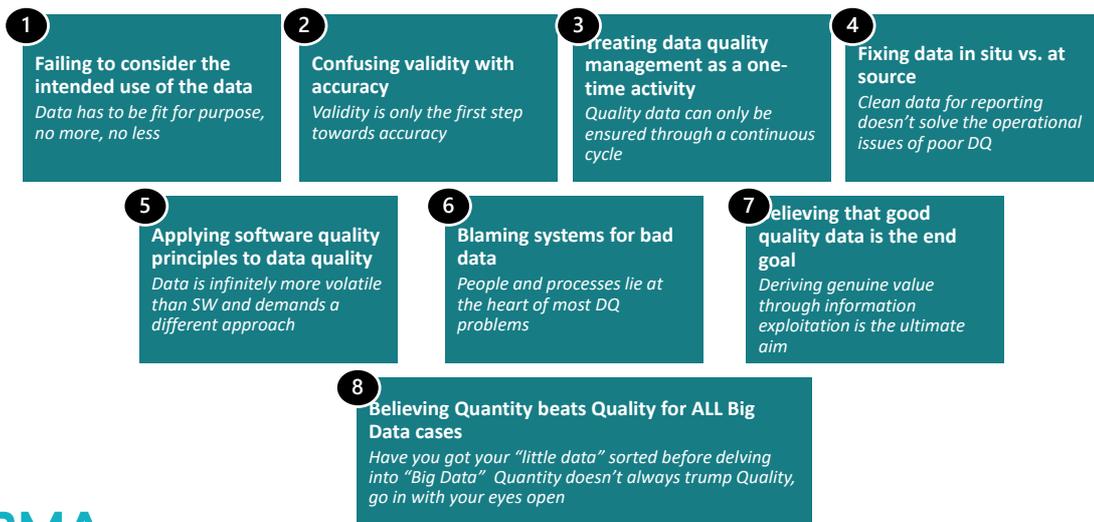
Controllability – make sure the remedial actions for each DQI are defined



T
A
R
M
A
C



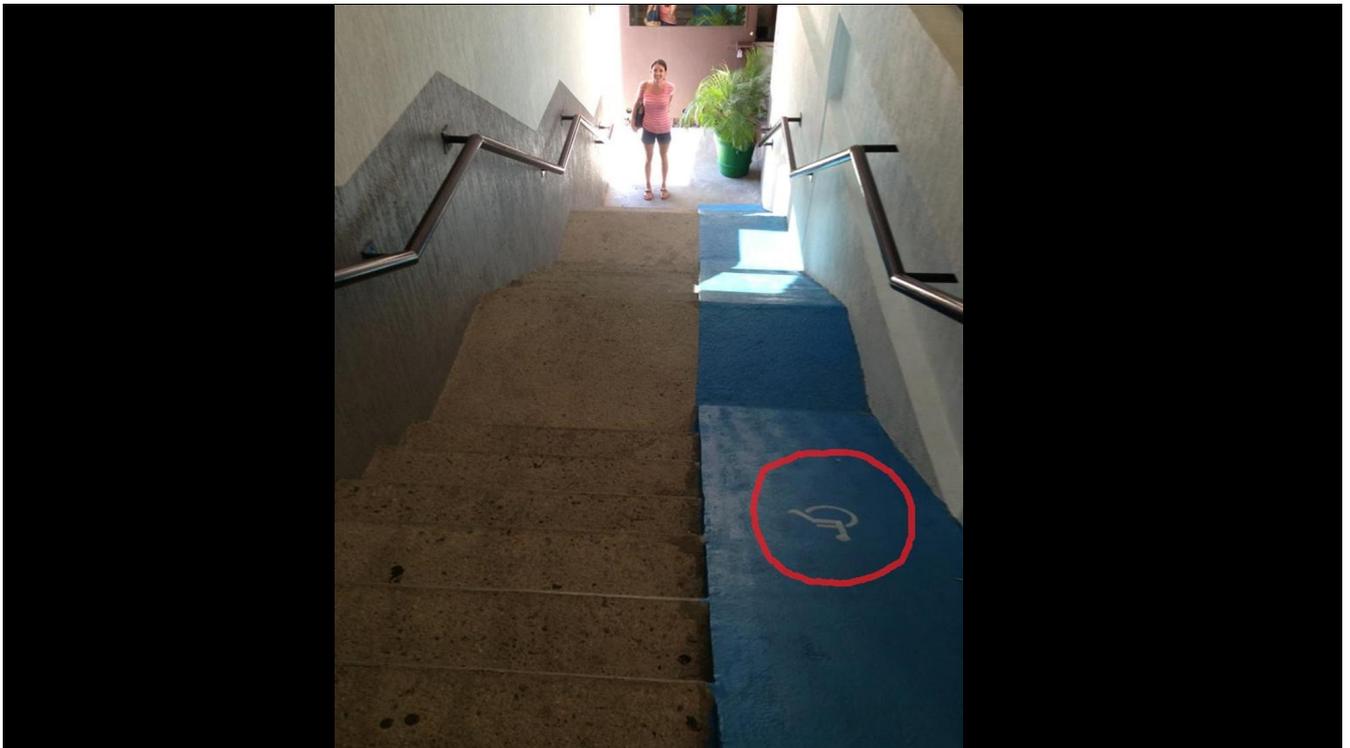
Common DQ Mistakes





...and finally ... It's not always obvious







Ref	Question	A	B	C	D	E
DQ1	When defining data quality indicators, care must be taken to ensure that they have	Measurability, Relevance & Acceptability	A direct link to the Data Governance strategy	Items in a dashboard showing their improvement over time	The core dimensions of Data Quality	Timeliness, Validity & Accuracy
DQ2	Which of these statements is true?	Data Quality Management is a synonym for Data Governance	Data Quality Management is a continuous process	Data Quality Management only addresses structured data	Data Quality Management is the application of technology to data problems	Data Quality Management is usually a one-off project
DQ3	The Data Quality Management cycle has four stages. Three are Plan, Monitor & Act. What is the fourth stage?	Improve	Prepare	Reiterate	Deploy	Manage
DQ4	Which of these is NOT a typical activity in Data Quality Management?	Defining business requirements & business rules	Analysing data quality	Creating inspection & monitoring processes	Identifying data problems & issues	Enterprise Data Modelling
DQ5	Which of these is NOT an expected role of a Data Quality Oversight Board?	Setting data quality improvement priorities	Establishing communications & feedback mechanisms	Data profiling & analysis	Producing certification & compliance policies	Approving data quality strategies
DQ6	According to DMBok, which of these is NOT a valid dimension of Data Quality?	Relevance	Timeliness	Currency	Completeness	Reasonableness
DQ7	Which of these is a key process in defining data quality business rules?	De-duplicating data records	Producing data management policies	Separating data that does not meet business needs from data that does	Matching data from different data sources	Producing data quality reports & dashboards
DQ8	A Data Quality Service Level Agreement (SLA) would normally include which of these?	Respective roles & responsibilities for data quality	A Business Case for data improvement	An enterprise data model	Detailed technical specifications for data transfer	A breakdown of the costs of data quality improvement
DQ9	'Top down' and 'bottom up' data analysis and profiling is best done together because	It balances business relevance and the actual state of the data.	It gets everyone involved.	It gives something for the architects to do while the profilers get on with the work	It allows the profiler to show the business the true state of the data	Data quality tools are more productive when they are effectively configured.

AFTER QUIZ 2

1. General (6)
2. Data Quality (9)

Maximum possible score = 15

60% (CDMP Associate) = 9

70% (CDMP Practitioner) = 11

80% (CDMP Master) = 12

112

Links to Additional Content

Additional Reading:

- *Out of the Crisis*, W. Edwards Deming
- *Data Quality, The Field Guide*, Dr. Thomas C Redman
- *Data Quality, The Accuracy Dimension*, Jack Olsen
- *Improving Data Warehouse and Business Information Quality*, Larry P. English
- *Memory Jogger 2*, Brassard et al. (contains useful summaries of quality management statistical tools)

Data Storage & Operations

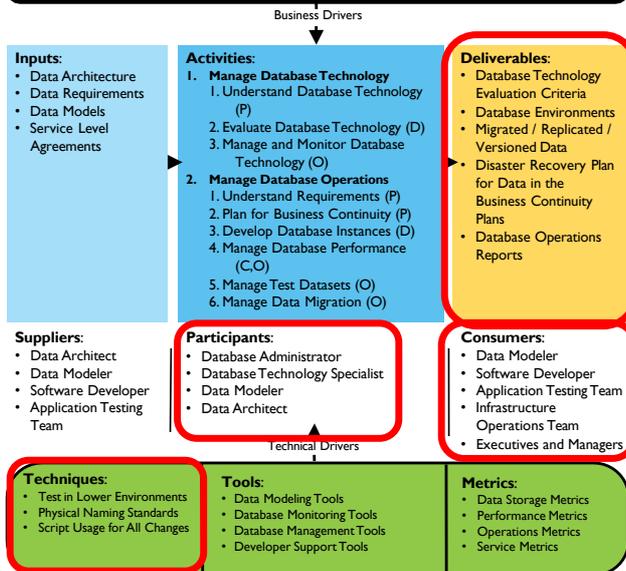
Data Management Process	2%
Big Data	2%
Data Architecture & Lifecycle Management	6%
Document, Records & Content Management	4-5%
Data Ethics	2%
Data Governance	11%
Data Integration & Interoperability	6-7%
Master & Reference Data Management	11%
Data Modelling & Design	11%
Data Quality	11%
Data Security	6%
Data Storage & Operations	4-5%
Data Warehousing & Business Intelligence	11%
Metadata Management	11%



Data Storage and Operations (DMBoK 2 revised)

Definition: The design, implementation, and support of stored data to maximize its value.

- Goals:**
1. Manage availability of data throughout the data lifecycle.
 2. Ensure the integrity of data assets.
 3. Manage performance of data transactions.





What is Data Operations Management?

"Data operations management is the development, maintenance, and support of structured data to maximize the value of the data resources to the enterprise. Data operations management includes two sub-functions: database support and data technology management" [DAMA DMBOK]

Key Points

Goals of Data Operations are based on:

- Protecting and ensuring the integrity of data assets
- Ensuring the availability of data throughout its lifecycle
- Optimize performance of database applications

Primary Deliverables

- Production database environments
- Controlled mechanisms and processes for implementation & changes to databases
- Mechanisms for ensuring availability, integrity and recoverability of data
- Mechanisms for detecting and reporting errors
- Mechanisms for ensuring performance according to service level agreements

Main Activities

Database Support

- Primarily undertaken by Production DBAs
- Ensuring the performance and reliability of the database, including performance tuning, monitoring, and error reporting.
- Implementing appropriate backup and recovery mechanisms.
- Implementing mechanisms for clustering and failover of the database (e.g. for 24/7 data availability).
- Implement mechanisms for data archiving.

		2-5 "9s" Availability			
		99%	99.9%	99.99%	99.999%
Data Technolo Managem	Daily	14m 24s	1m 26s	8s	1s
	Weekly	1h 40m 48s	10m 4s	1m 0s	6s
	Monthly	7h 18m 17s	43m 49s	4m 22s	26s
	Quarterly	21h 54m 52s	2h 11m 29s	13m 8s	1m 18s
	Yearly	3d 15h 39m 29s	8h 45m 56s	52m 35s	5m 15s



	Hot Data		Warm Data			Cold Data	
	Amazon ElastiCache	Amazon DynamoDB	Amazon Aurora	Amazon Elasticsearch	Amazon EMR (HDFS)	Amazon S3	Amazon Glacier
Average latency	ms	ms	ms, sec	ms,sec	sec,min,hrs	ms,sec,min (~ size)	hrs
Data volume	GB	GB-TBs (no limit)	GB-TB (64 TB Max)	GB-TB	GB-PB (~nodes)	MB-PB (no limit)	GB-PB (no limit)
Item size	B-KB	KB (400 KB max)	KB (64 KB)	KB (1 MB max)	MB-GB	KB-GB (5 TB max)	GB (40 TB max)
Request rate	High - Very High	Very High (no limit)	High	High	Low - Very High	Low - Very High (no limit)	Very Low
Storage cost GB/month	\$\$	cc	cc	cc	c	c	c/10
Durability	Low - Moderate	Very High	Very High	High	High	Very High	Very High
	Hot Data		Warm Data			Cold Data	



Factors affecting Availability vs Performance

Manageability

The ability to create and maintain an effective environment.

Recoverability

The ability to re-establish service after interruption, and correct errors caused by unforeseen events or component failures.

Reliability

The ability to deliver service at specified levels for a stated period.

Serviceability

The ability to determine the existence of problems, diagnose their cause and repair/solve the problems.

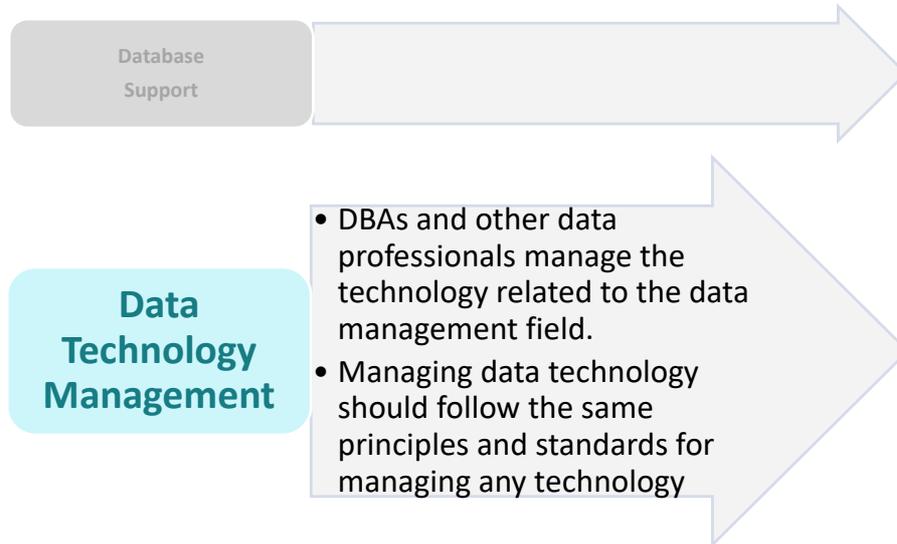


Causes of poor database performance

- Memory allocation (buffer/cache for data)
- Locking and blocking
- Failure to update database statistics (SysTables)
- Poor SQL coding
- Poor design of Views (yielding the query from hell)
- Inappropriate indexing
- Application activity
- Increase in the number, size or use of databases
- Database volatility



Main Activities



Technology Architecture Components - Considerations

Current

Products currently supported and used.

Deployment Period

Products to be deployed for use in the next 1-2 years.

Strategic Period

Products expected to be available for use in the next 2+ years.

Retirement

Products the organisation has retired or intends to retire this year.

Preferred

Products preferred for use by most applications.

Containment

Products limited to use by certain applications.

Emerging

Products being researched and piloted for possible future deployment.

Ref	Question	A	B	C	D	E
DSO 1	Which of the following activities are performed by data operations staff?	Implement and control database environments, plan for data retention, keep track of database licenses, monitor and tune database performance	Grant access to tables, rewrite SQL statements	Clean data that is of bad quality	Manage the tape libraries	Tune the file systems
DSO 2	The goals of data operations include which of the following?	Assuring the quality of the structured data assets, taking backups and managing security of the database	Assuring availability of the data throughout its lifecycle, protection and integrity assurance of structured data assets and performance optimization of database transactions	Assuring the performance of the network and storage devices, the quality of the SQL statements and the selection of DBMS platform	Assuring backups are taken, managing the performance of SQL and checking data quality	Providing the right database access rights, solving software bugs and managing database logs
DSO 3	The data operations team assures that the data is recoverable by...	Making sure the disks are checked regularly for write errors.	Guaranteeing the applications take proper exports of the data.	Defining and executing the data recovery plan.	Maintaining a test, development and production environment.	Analysing database error logs.



AFTER QUIZ 3

1. General (6)
2. Data Quality (9)
3. Data Storage & Operations (3)

Maximum possible score = 18

60% (CDMP Associate) = 11

70% (CDMP Practitioner) = 13

80% (CDMP Master) = 15

Reference & Master Data Management

Data Management Process	2%
Big Data	2%
Data Architecture & Lifecycle Management	6%
Document, Records & Content Management	4-5%
Data Ethics	2%
Data Governance	11%
Data Integration & Interoperability	6-7%
Master & Reference Data Management	11%
Data Modelling & Design	11%
Data Quality	11%
Data Security	6%
Data Storage & Operations	4-5%
Data Warehousing & Business Intelligence	11%
Metadata Management	11%



Reference and Master Data (DMBoK 2 revised)

Definition: Managing reconciled and integrated data through stewardship and semantic consistency in support of enterprise-wide needs to share its data assets.

Goals:

1. Ensuring the organization has complete, consistent, current, authoritative Master and Reference Data across organizational processes.
2. Enabling Master and Reference Data to be shared across enterprise functions and applications.
3. Lowering the cost and reducing the complexity of data usage and integration through standards, common data models, and integration patterns.

Business Drivers

Inputs:

- Candidate Data Stores and Values
- Cross Functional Requirements
- Industry Standards
- Metadata
- Purchased Data and/or Open Data and Code Sets
- Business Rules

Activities:

1. Define Drivers and Requirements (P)
2. Evaluate and Assess Data Sources (P)
3. Define Architectural Approach (D)
4. Model Data (D)
5. Define Stewardship and Maintenance Processes (C)
6. Establish Governance Policies (C)
7. Implement Data Sharing / Integration Services (D,O)
 1. Acquire Data Sources for Sharing
 2. Publish Reference and Master Data

Deliverables:

- Master and Reference Data Requirements
- Data Models and Integration Patterns
- Reliable Reference and Master Data
- Reusable Data Services
- Data Exception Reports

Suppliers:

- Subject Matter Experts
- Data Stewards
- Application Developers
- Data Providers
- Business Analysts
- Infrastructure Systems Analysts

Participants:

- Data Analysts
- Data Modelers
- Data Stewards
- Data Integrators
- Data Architects
- Data Quality Analysts

Consumers:

- Subject Matter Experts
- Data Integrators
- Application Users
- Application Developers
- Solution Architects

Technical Drivers

Techniques:

- Conditions-of-use agreements
- Business key cross references
- Processing Log analysis

Tools:

- Data Integration Tools
- Data Remediation Tools
- Operational Data Stores
- Data Sharing Hubs
- Data Modeling Tools
- Metadata Repositories
- MDM Application Platforms

Metrics:

- Data Quality and Compliance
- Data Change Activity
- Data Ingestion and Consumption
- Data Sharing Availability
- Data Steward Coverage
- Data Sharing Volume and Usage



What is Event / Transaction Data?

EVENT DATA EXAMPLE:

“Bob bought a Cadburys Dairy Milk bar from Morrison's on Monday 3rd Jan at 4pm and paid using cash.”



WHO	WHAT	WHERE	WHEN	HOW	QUANTITY	AMOUNT
Bob Smith	Dairy Milk bar	Morrison's, Bath	16:00 Monday 3 rd January 2022	Cash	1	£7.60

CUSTOMER CODE	PRODUCT CODE	VENDOR CODE	DATE	PAYMENT METHOD	QUANTITY	AMOUNT
BS005	CONF101	WMBATH	2022-01-03 16:00	CASH	1	£7.60



About Event Data

AKA Transaction data

Describes an action (a verb)
E.g., “buy”

Will include measurements about the action:

- Quantity bought
- Amount paid

Includes information:

Identifying the nouns that were involved in the event (the Who / What / Where / When / How and maybe even the Why):

- Bob Smith
- Dairy Milk bar
- Morrisons, Bath
- 16:00 Monday 3rd Jan 2022
- Cash

Does not include information:

Describing the nouns:

- Roberta is female, aged 25 and works for British Airways, drives an Audi,
- That Dairy Milk is a special offer, extra-large, 850g jumbo bar, half milk, half dark, the recipe is ...
- The address of Morrisons Bath is York Place, London Road, Bath, BA1 6AE, the manager is Gary Smith, ...
- Monday 3rd Jan 2022 is a public holiday
- Other methods of payment include store card, debit card, Apple pay ...



What is...

MASTER & REFERENCE DATA?

- Defines and describes the nouns (things) of the business.
e.g., Customer, Risk, Branch, Loan, Field, Well, Rig, Product, Store, Therapeutic Area, Adverse Event, etc.
- Data about the “things” that will participate in events.
- Provides *contextual information* about events / transactions.
- This Data is stored in **many** systems
 - Application package systems
 - Line of Business Systems
 - Spreadsheets
 - SharePoint Lists,



MASTER DATA MANAGEMENT (MDM)?

- The ongoing reconciliation and maintenance of master data.
- Control over master data values to enable consistent, shared, contextual use across systems, of the most accurate, timely, and relevant version of truth about essential business entities.
(DAMA)
- Comprises a set of processes and tools that consistently defines and manages the non-transactional data entities of an organisation.
(Wikipedia)

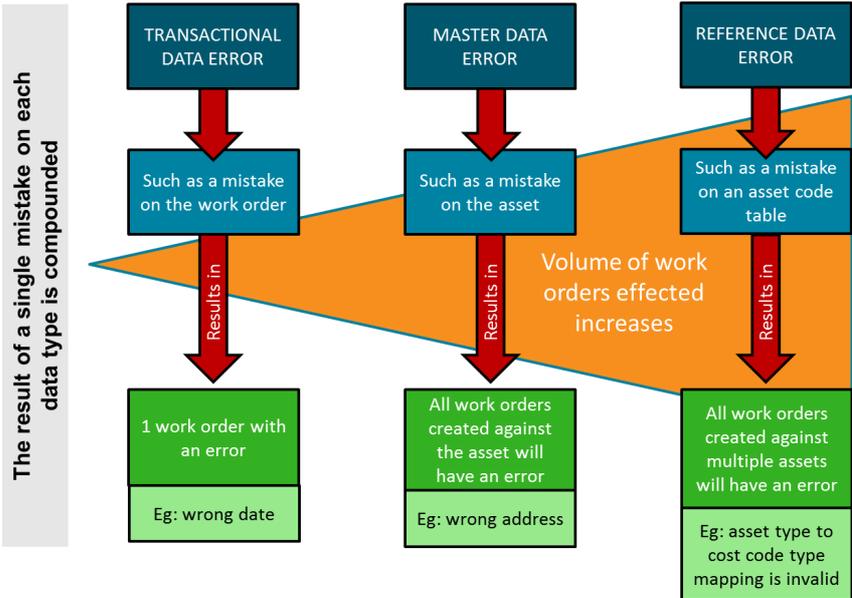
Reference vs Master Data

Characteristic	Reference Data	Master Data
Number of Values	Low, Fixed/Known	Medium-High, Variable/unknown
Volatility	Very low rate of change	Medium-Frequent rate of change
Source	External (mostly); Standards bodies	Internal (mostly)
Number of Attributes per data area	Very low; typically code type/value/description	High
Federation/ownership of attributes per business area	None (all attributes)	Much federation: e.g.. Business Area A masters Attribute 1/Attribute 2, Business Area B masters Attributes 3/4/5
Ease of Governance	Easy	Harder, Higher number of stakeholders
Tool Complexity & Cost	Low (Ref data only) Code (s/w) distribution (e.g. dropdown lists)	High (MDM tools often have ref-data management also)

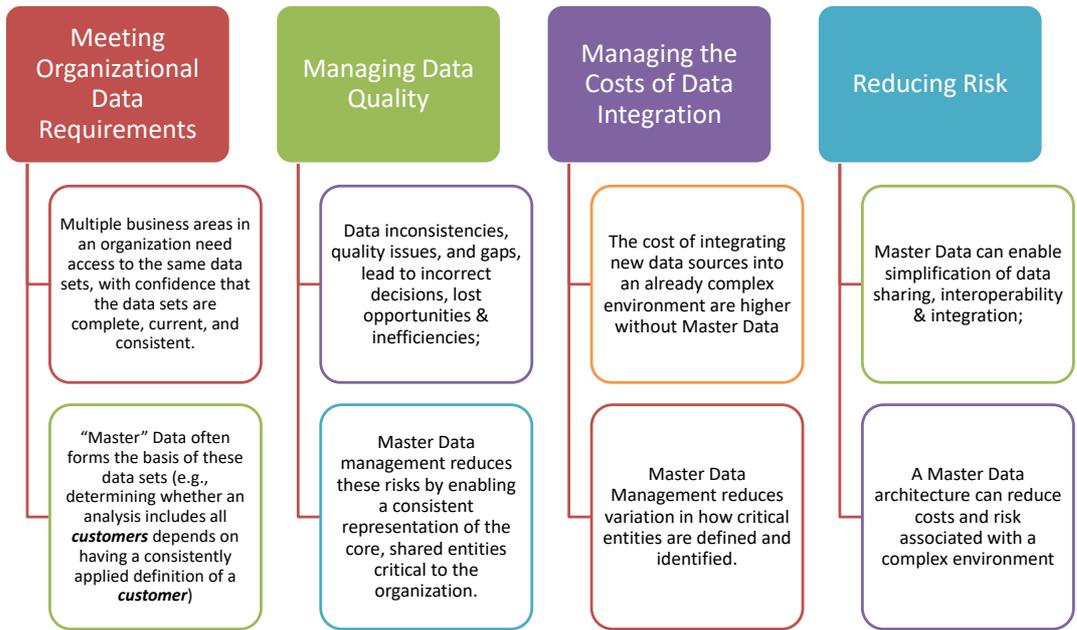




Importance of Reference Data

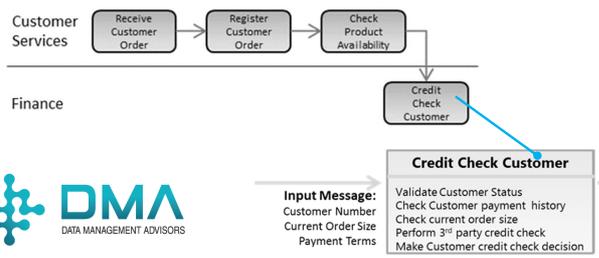
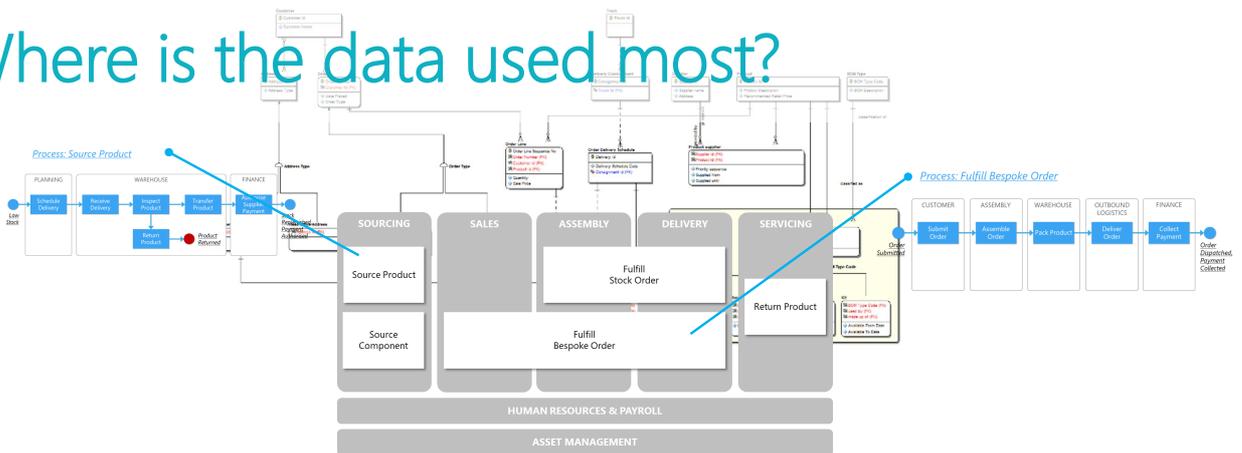


MDM Business Drivers





Where is the data used most?



Business Processes

	Product development	Planning & Sales	Industrial engineering	Order management	Manufacturing	Logistics	Marketing
Product	C	R	U	U	U		
Product Part	C	R	R	U	U		
Manufacturing Plant	U		C	R	R	U	
Customer	R	C		U	R	U	U
Sales Item	C	C	C	U		U	U
Assembly Structure	U		C		U		
Sales Order		U		R	U	U	U
Production Order			U	C	U	U	U
Industrial Product					C	R	U
Shipping						C	
Customer's Interest		U				P / 138	C

Major Entities / Data Subject Areas

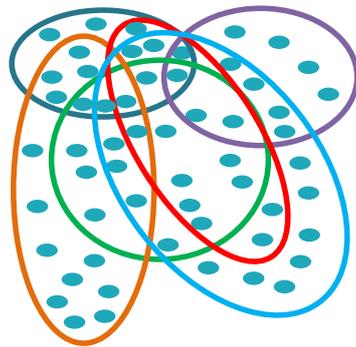
Process and Data ARE related

Create
 Read
 Update
 Delete

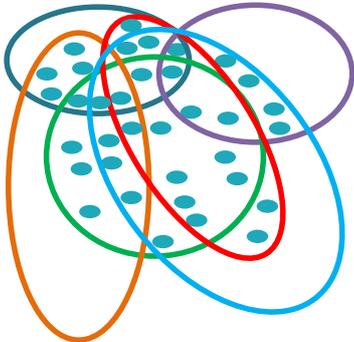
Major Entities / Data Subject Areas	Business Processes						
	Product development	Marketing & Sales	Industrial preparation	Order management	Manufacturing	Logistics	Invoicing
Product	C	R	U	U	U		
Product Part	C	R	R	U	U		
Manufacturing Plant	U		C	R	R	U	
Customer	R	C		U	R	U	U
Sales Item	C	C	C	U		U	U
Assembly Structure	U		C		U		
Sales Order		U		R	U	U	U
Production Order			U	C	U	U	U
Individual Product					C	R	U
Shipping						C	
Customer's Invoice		U					C



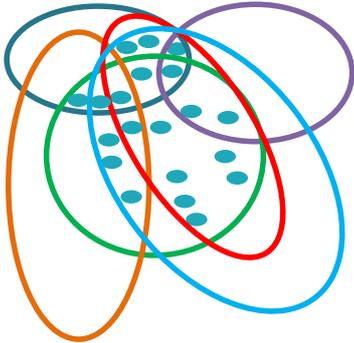
What should be mastered?



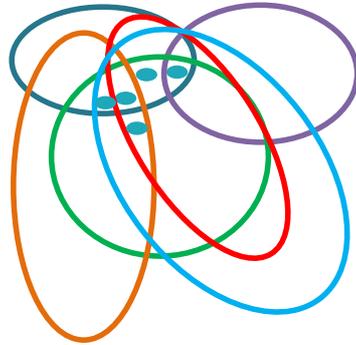
What should be mastered?



What should be mastered?



What should be mastered?



MDM architectures

Standard
"Hub"
architectures

1. REPOSITORY

2. REGISTRY

A key difference is the number of fields that are stored centrally

3. HYBRID

4. VIRTUALISED

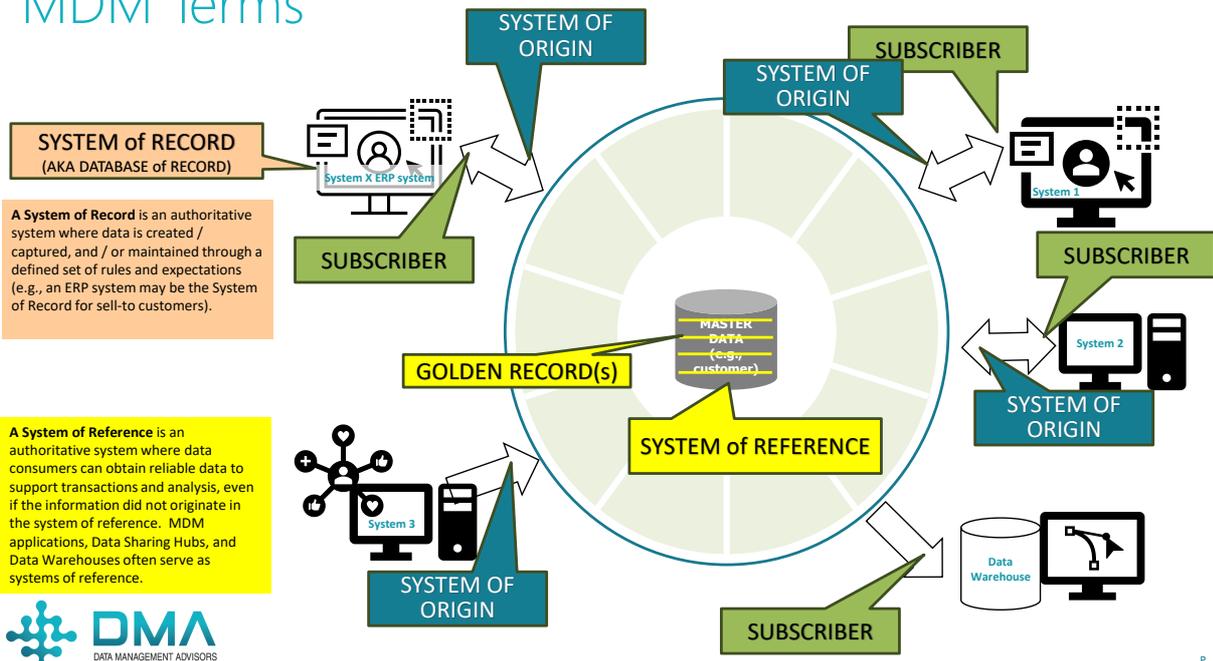
Example: Customer

Customer code	First name	Last name	Date of birth	Preferred delivery address post code	Preferred delivery address line 1	Credit rating	Occupation	Car
JBS005	Roberta	Smith	1985-12-25	BA1 7LA	Royal Crescent	A	Information Architect	Audi R8

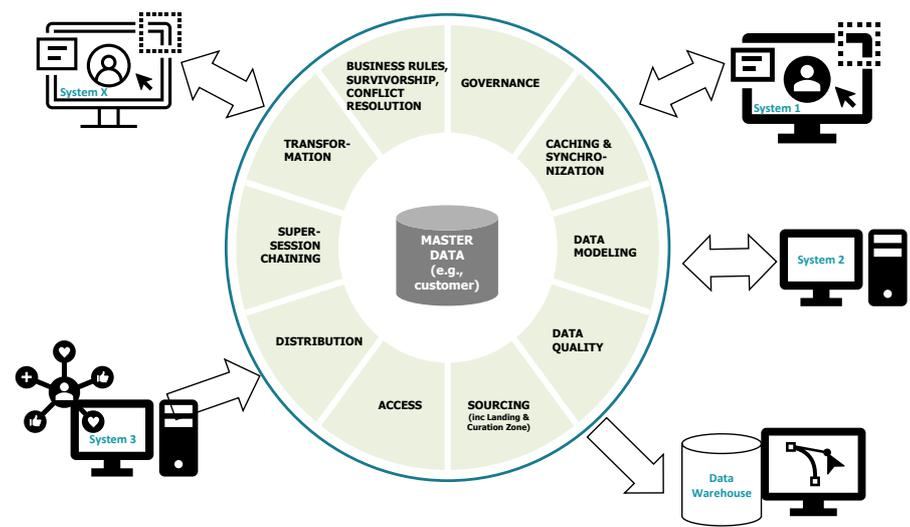
Example: Customer

Customer code	First name	Last name	Date of birth	Preferred delivery address post code	Preferred delivery address line 1	Credit rating	Occupation	Car
BS005	Roberta	Smith	1985-12-25	BA1 7LA	Royal Crescent	A	Information Architect	Audi R8	

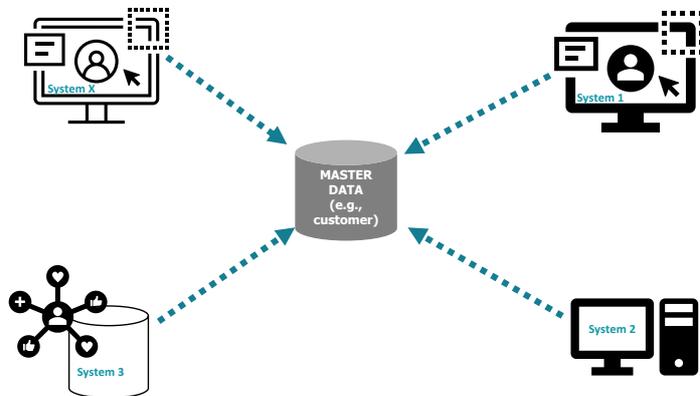
MDM Terms



Typical MDM Components



MDM Architecture Implementation Styles REGISTRY



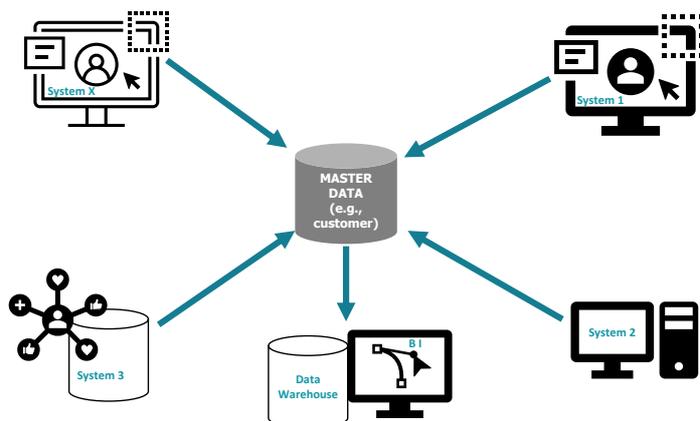
- Low control, autonomous environments
- The most difficult MDM style to implement data governance.
- Non intrusive to edge applications
- Emphasis is on remote data and Application to Application integration
- Distributed governance
- Faster to implement than coexistence & centralised

BENEFITS

- Can be used to analyse the data while avoiding the risk of overwriting information in the source systems.
- Helps avoid potential compliance failure or other regulatory repercussions
- Provides a read-only view of data without modifying master data and is a useful way to remove duplications and gain consistent access to your master data
- Low-cost, rapid data integration with the benefit of minimal intrusion into your application systems.

P / 149

MDM Architecture Implementation Styles CONSOLIDATION



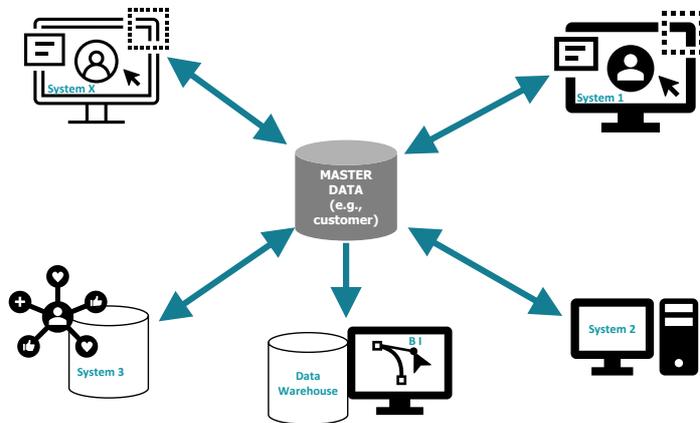
- Ideal for reporting or analytics that reside in a BI / DW environment
- Non intrusive to the business
- BI is the “main” use
- Benefits primarily BI use
- No cleaning of source data
- Higher data latency

BENEFITS

- Pulls master data from several existing systems into a single managed Master Data Management hub.
- Data is cleansed, matched and integrated to offer a complete single record (Golden Record) for one or more master data domains.
- Inexpensive and quick to set up, providing a fast and efficient way to facilitate enterprise-wide reporting.
- Mainly used for analysis, providing a trusted source of data for reporting and analytics.

P / 150

MDM Architecture Implementation Styles COEXISTENCE



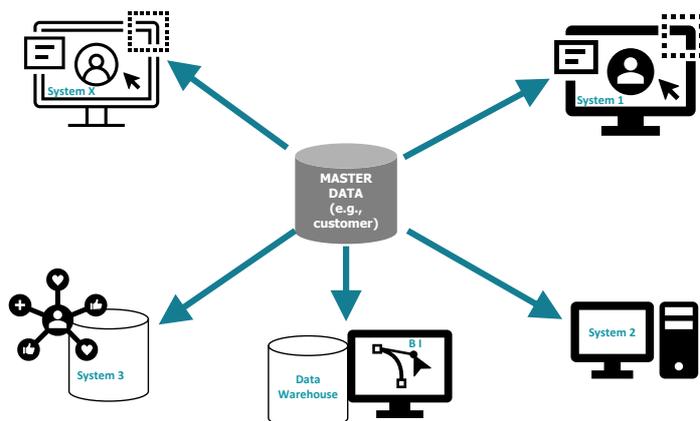
- Large scale distributed model
- Greatest change to information infrastructure
- Greatest need to mirror data
- Global AND local governance
- Risk over control & security
- Focused on shared services
- Lives the “single version of the truth”

BENEFITS

- Data is mastered in source systems and then synchronized with the hub.
- Data can coexist harmoniously and still offer a single version of the truth.
- Quality of master data is improved, and access is faster.
- Reporting is easier as all master data attributes are in a single place.
- Consolidation style can naturally evolve into a Coexistence style if business decides it requires the advantage of linking centrally governed, quality data back to the source systems.

P / 151

MDM Architecture Implementation Styles CENTRALISED

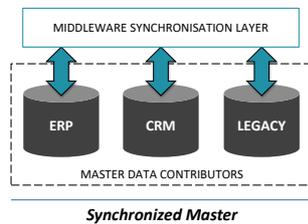


- High control top-down environments.
- Centralised governance - the easiest MDM style to implement data governance based on controls placed on persistent data.
- Greatest change to application infrastructure
- Highly invasive to the business
- Greatest control over access & security
- Focused on common services
- Lives the “single version of the truth”

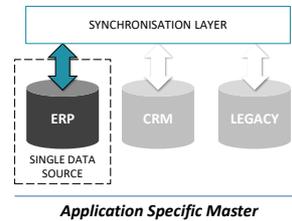
BENEFITS

- Master data is accurate and complete at all times.
- Security and visibility policies at a data attribute level can be supported.
- Centralized set of master data for one or more domains.
- MDM now is the system of origin, can leverage powerful data governance capabilities of the MDM.
- Centralized style can evolve from Consolidation or Coexistence styles. P / 152

MDM... A Hub is not the only way



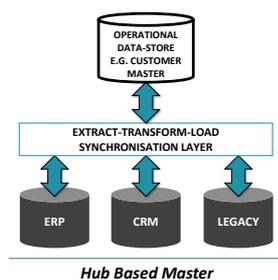
- Multiple operational systems acting as master data contributors
- Real-time information availability
- Well suited to enterprises where data is stored across multiple source systems
- Well suited to low data velocity operations
- **Like Coexistence**



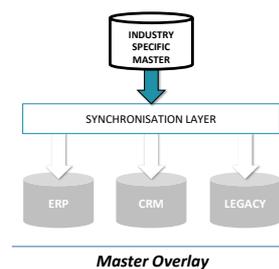
- One operational system as master data provider
- Well suited to enterprises where data is primarily stored in a single source system
- Support from many enterprise vendors
- **Hybrid, partially like Coexistence**



MDM... A Hub is not the only way



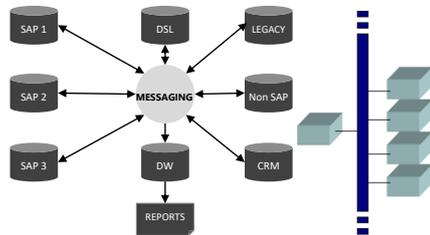
- Operational Hub structure that overlaps operational and analytical environments
- Supports concept of an Enterprise Data Warehouse
- Multiple systems acting as data providers
- Appropriate for low data latency & velocity operations
- Requires careful data quality management
- **Coexistence**



- Stand-alone, application-neutral master data
- Industry-specific data model
- Well suited to vertical industries in aligning front-office & back-office systems in real time
- Now most often seen for **Reference Data**
- **Centralised**

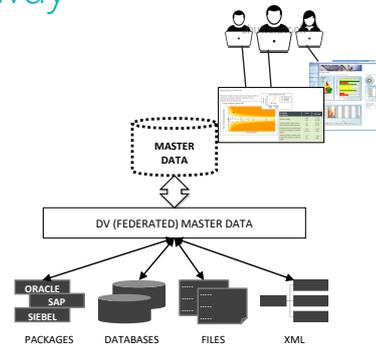


MDM... A Hub is not the only way



Real Time Data Movement

- Data is moved between the various systems using a messaging integration hub or ESB
- Data movement is done in real time
- Movement of data can be one or two way as required
- **Coexistence via a message-based architecture**



Data Virtualization (Federated) Master Data

- Virtual MDM hub created via DV layer
- DV layer informed from Logical Data Model
- Powerful data integration to access multiple disparate systems
- Rapid time to solution
- Outstanding prototyping, proof of value approach
- Compliments other Data Integration approaches
- **Could be Consolidated, Coexistence or Centralised**

Master Data Match Rules

Three primary scenarios:

Rules around the matching, merging and linking of data from multiple systems about the same person, group, place or thing.

1. Duplicate identification match rules

Focus on a specific set of fields that uniquely identify an entity and identify merge opportunities without taking automatic action. Business Data stewards can review these occurrences and decide to act on a case-by-case basis.

Less used

2. Match-merge rules

Match records and merge the data from these records into a single, unified, reconciled and comprehensive record. If the rules apply across data sources, create a single unique and comprehensive record in each database.

Most commonly used

3. Match-link rules

Identify and cross-reference records that appear to relate to a master record without updating the content of the cross-referenced record. Match-link rules are easier to implement and much easier to reverse.

Used extensively in Police / security services, local authorities



.. and Data Enrichment can also help with business entity resolution !

Master Data Matching

Deterministic matching algorithm

- Matches exact character to character of one or more fields (**Exact string match**)
- Has a discrete all or nothing outcome
- All identifiers being matched have equal weight
- Is better suited when there is no great consequence to an error in matching

Probabilistic matching algorithm

- Each variable to be matched is assigned a weight based on its discriminating power
- A score is assigned based on weight and degree of match
- Individual attribute matching scores are used to create a match probability percentage
- Following the matching process there are typically records requiring manual review and decisioning

MDM matching algorithms benefit from the data characteristics of:

- High validity of the data
- Distinctiveness across the population of data
- High level of comparability of the data elements
- Structural heterogeneity of data elements



Master Data Matching



A **TRUE NEGATIVE**

WHEN 2(OR MORE) RECORDS ARE **NOT** MATCHED WHEN THEY **SHOULD NOT** HAVE BEEN MATCHED

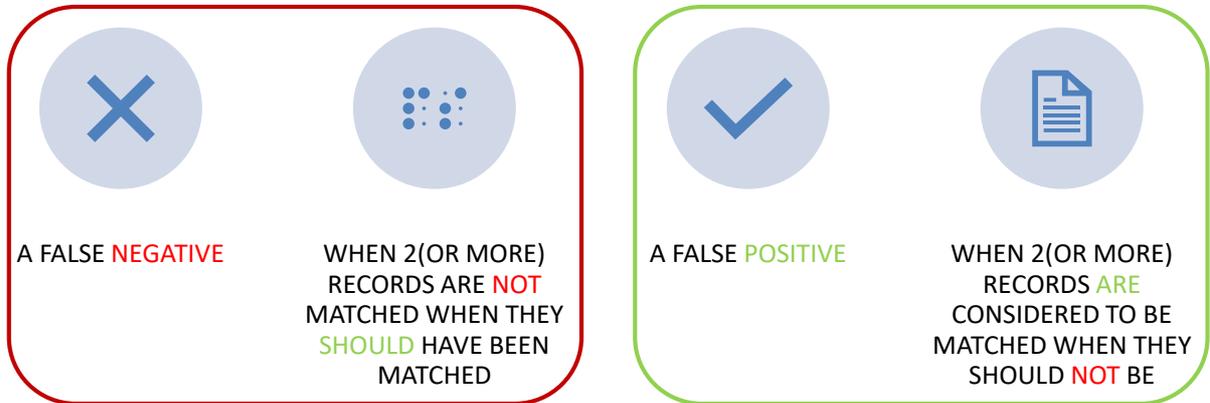


A **TRUE POSITIVE**

WHEN 2(OR MORE) RECORDS **ARE** CONSIDERED TO BE MATCHED WHEN THEY **SHOULD BE**



Master Data Matching



Single Domain & Multi Domain MDM

SINGLE DOMAIN

- E.g. Product (PIM), Customer (UCM), Vendor, Laboratory (LIM)
- Rich Data Model (frequently an Industry data model)
- Incorporate specific “domain” features ...
 - Rich functionality & logic for the specific data domain
 - House holding
 - Address chaining
 - Company Hierarchies
 - Tailored data matching (deterministic or probabilistic)
- Engineered to exploit 3rd party data sources
 - GB PAF, USPS Zip+4, Electoral Roll, CCJ ...
 - D&B, Liquidations, Companies House, ...
 - CAS, COSHH, ...
- Interfaces to domain LoB apps

MULTI DOMAIN

- Highly configurable MDM solutions
- Informed by your Data Model
- Fewer specific data domain features
- Standard MDM processes
- Interfaces to mainstream apps
- Results in fewer MDM solutions throughout the Enterprise

Implementation: Operational vs Analytical MDM

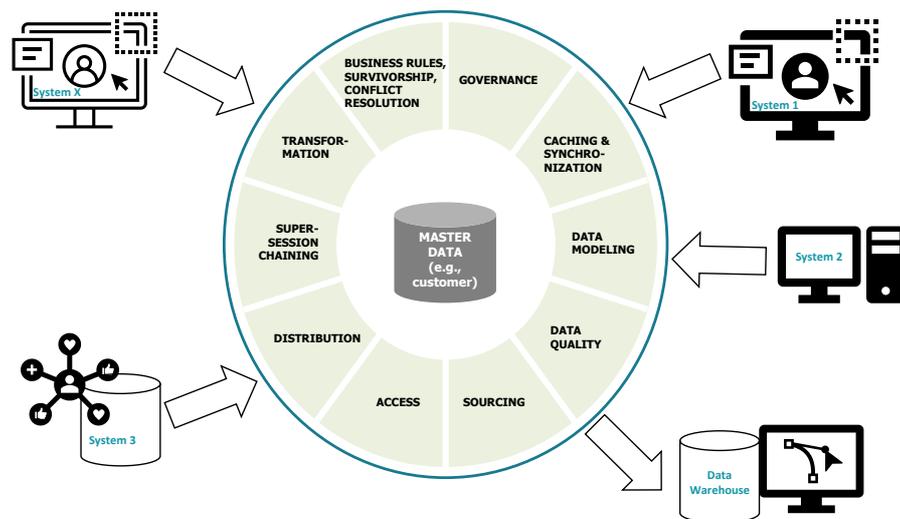
ANALYTICAL

- MDM (usually a hub) created for use in Data Warehouse / Business Intelligence solutions
- Data integration (e.g., ETL) from source systems still required
- Data cleanse, consolidate, merge, de-duplicate, enrich etc. still required
- Master Data used for analytical purposes to ensure consistency of analyses, findings etc.
- Master Data NOT used by live operational line of business solutions, therefore operational systems do not see any benefits of MDM approach
- NOT essential to address business rules (such as survivorship, channels, governance) for operational MDM
- Easier (organizationally) to implement
- Impact of failure only on BI activities



MDM Implementation

Analytical



Implementation: Operational vs Analytical MDM

OPERATIONAL

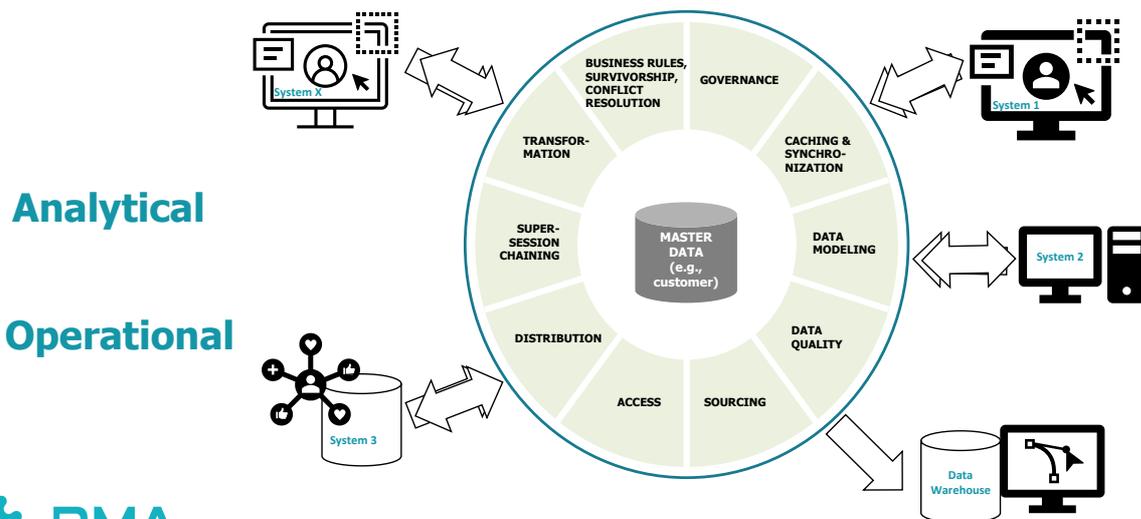
- MDM (whatever is most appropriate architecture) created for use in LIVE operational Business solutions
- Data integration (e.g., ETL) from source systems essential
- Data cleanse, consolidate, merge, de-duplicate, enrich etc. essential
- Master Data is used by the live operational line of business solutions
- Essential to address business rules (such as survivorship, conflict resolution, channels, attribute federation, data governance) for operational MDM
- More difficult (technically & organizationally) to implement
- Impact of failure severe
- Benefits extensive across the organization

ANALYTICAL

- MDM (usually a hub) created for use in Data Warehouse / Business Intelligence solutions
- Data integration (e.g., ETL) from source systems still required
- Data cleanse, consolidate, merge, de-duplicate, enrich etc. still required
- Master Data used for analytical purposes to ensure consistency of analyses, findings etc.
- Master Data NOT used by live operational line of business solutions, therefore operational systems do not see any benefits of MDM approach
- NOT essential to address business rules (such as survivorship, channels, governance) for operational MDM
- Easier (organizationally) to implement
- Impact of failure only on BI activities



MDM Implementation



Analytical

Operational



It's NOT "The Field of Dreams"

- Projects cannot afford to wait for Master Data to be delivered
- Subsets of the Master Data must be delivered "just in time" as the Business projects need to use them
- This is NOT "The Field Of Dreams"
- This means **aligning the MDM initiatives with the Business Projects**



P / 165

Conclusions

Don't believe you can "build it and they will come"

An MDM initiative needs to be just ahead of the business projects & deliver just in time aligned to the business initiatives.

Data models for MDM programs are essential

Determine data areas to master.

Configure MDM product.

Governance .. E.g. extended with rich metadata including Ownership, Source systems, Transformations etc.

Data Governance is vital

Operational MDM initiatives **cannot** be successful without addressing Data Governance for that data subject area.

Beware the hype on MDM technologies

A hub is NOT the only way.

Rarely will one MDM technology be appropriate / best for ALL Data Subject Areas.



P / 166



CHOOSE WISELY

Key Points (DMBoK 2 Rev)

- Shared Reference and Master data **belongs to the organization**, not to a particular application or department
- Reference and Master data management is an **on-going data quality improvement program**; its goals cannot be achieved by one project alone
- **Business data stewards** are the authorities **accountable** for controlling **reference data** values. Business data stewards work with data professionals to improve the quality of reference and master data
- Request, communicate, and, in some cases, approve of changes to reference data values before implementation
- Golden data values (**Golden records**) represent the organization's **best efforts** at determining the most accurate, current, and relevant data values for contextual use. New data may prove earlier assumptions to be false. Therefore, apply matching rules with caution, and ensure that any changes that are made are reversible
- Replicate master data values only from the **system of record**
- *A **System of Record** is an authoritative system where data is created/captured, and/or maintained through a defined set of rules and expectations (e.g., an ERP system may be the System of Record for sell-to customers).*
- *A **System of Reference** is an authoritative system where data consumers can obtain reliable data to support transactions and analysis, even if the information did not originate in the system of reference. MDM applications, Data Sharing Hubs, and Data Warehouses often serve as systems of reference.*

Links to Additional Content

- Berson, Alex and Larry Dubov. Master Data Management and Customer Data Integration for a Global Enterprise. McGraw-Hill, 2007. ISBN 0-072-26349-0. 400 pages.
- Brackett, Michael. Data Sharing Using A Common Data Architecture. New York: John Wiley & Sons, 1994. ISBN 0-471-30993-1. 478 pages.
- Chisholm, Malcolm. Managing Reference Data in Enterprise Databases: Binding Corporate Data to the Wider World. Morgan Kaufmann, 2000. ISBN 1-558-60697-1. 389 pages.
- Dreibelbis, Allen, Eberhard Hechler, Ivan Milman, Martin Oberhofer, Paul van Run, and Dan Wolfson. Enterprise Master Data Management: An SOA Approach to Managing Core Information. IBM Press, 2008. ISBN 978-0-13-236625-0. 617 pages.
- Dyche, Jill and Evan Levy. Customer Data Integration: Reaching a Single Version of the Truth. John Wiley & Sons, 2006. ISBN 0-471-91697-8. 320 pages.
- Finkelstein, Clive. Enterprise Architecture for Integration: Rapid Delivery Methods and Techniques. Artech House Mobile Communications Library, 2006. ISBN 1-580-53713-8. 546 pages.
- Loshin, David. Master Data Management. Morgan Kaufmann, 2008. ISBN 98-0-12374225-4. 274 pages.



Ref	Question	A	B	C	D	E
MDM1	What is a common motivation for Reference & Master Data Management?	The need to improve data quality and data integrity across multiple data sources	The need to build a Data Dictionary of all core data entities & attributes	Regulatory acts such as BCBS239, GDPR and SOX	The need to consolidate all data into one physical database	Business Intelligence & Data Warehousing
MDM2	Which of these is a valid definition of Master Data?	Data that if missing or incorrect will cause transactions and processes to fail	Data that is only held in one data source	Data that other data sits hierarchically beneath	Data that rarely, if ever, changes	Data about the business entities that provide context for business transactions
MDM3	Which of these is a valid definition of Reference Data?	Data that is fixed and never changes	Data used to classify or categorize other data	Data that provides metadata about other data entities	Data that is widely accessed and referenced across an organisation	Data that has a common and widely understood data definition
MDM4	Which of the following is NOT a primary Master Data Management area of focus?	Generating a golden record / best version of the truth	Identifying duplicate records	Producing read only versions of key data items	Providing access to golden data records	Producing clear data definitions for Master Data
MDM5	A strong argument for pursuing a Reference Data and/or Master Data management initiative is:	It will not require a lot of time or effort	They are essential functions in the data management framework	Job security for the data people	By centralizing the management of Reference and Master data, the organization can conform critical data needed for analysis	Application retirement
MDM6	A common driver for initiating a Reference Data Management program is:	It will improve data quality and facilitate analysis across the organization	It can be a one-time-only project	Managing codes and descriptions requires little effort and low cost	It will consolidate the process of securing third party code sets	Application simplification
MDM7	Reference Data Management includes defining relationships within and across domain value lists.	TRUE	FALSE			
MDM8	Which one of the following statements is true?	Business data stewards maintain lists of valid data values for master data instances.	Managing reference data requires the same activities and techniques as does managing master data.	Reference Data Management involves identifying the 'best' or 'golden' record for each domain.	Master Data Management requires techniques for splitting or merging an instance of a business entity.	Operational Master Data Management can be introduced before the Data is Governed
MDM9	An authoritative system where data is created/captured, and/or maintained through a defined set of rules and expectations is called	A System of Record.	A System of Origin.	A System of Referential Integrity.	A System of Retirement.	A System of Systems.

AFTER QUIZ 4

1. General (6)
2. Data Quality (9)
3. Data Storage & Operations (3)
4. Master & Reference Data (9)

Maximum possible score = 27

60% (CDMP Associate) = 16

70% (CDMP Practitioner) = 19

80% (CDMP Master) = 22

172

Data Warehousing & Business Intelligence

Data Management Process	2%
Big Data	2%
Data Architecture & Lifecycle Management	6%
Document, Records & Content Management	4-5%
Data Ethics	2%
Data Governance	11%
Data Integration & Interoperability	6-7%
Master & Reference Data Management	11%
Data Modelling & Design	11%
Data Quality	11%
Data Security	6%
Data Storage & Operations	4-5%
Data Warehousing & Business Intelligence	11%
Metadata Management	11%

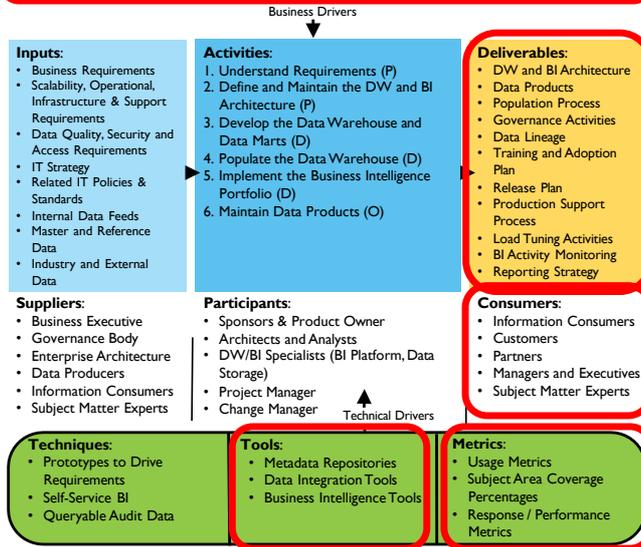


Data Warehousing and Business Intelligence (DMBoK 2 revised)

Definition: Planning, implementation, and managing an integrated data system to support knowledge workers engaged in reporting, query, and analysis.

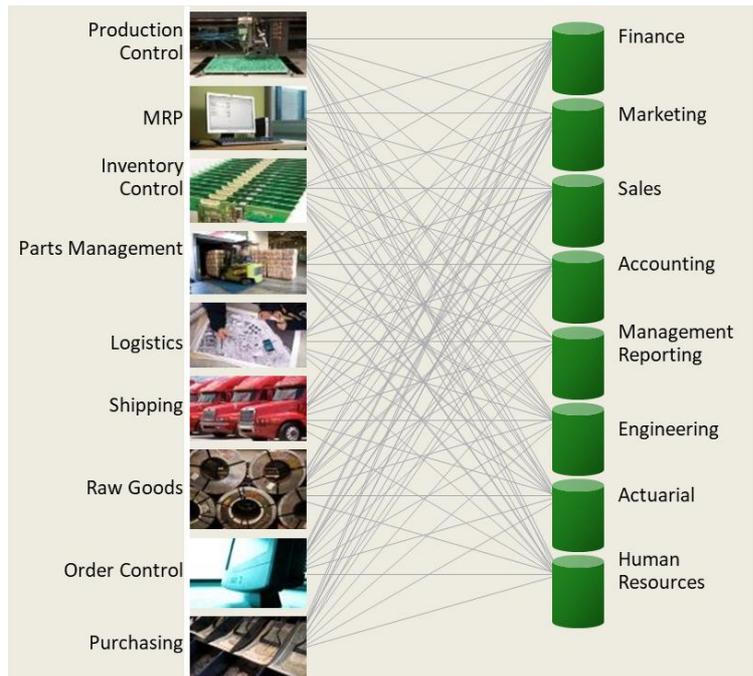
Goals:

- To build and maintain the data system with the technical and business requirements needed to deliver integrated data that supports operational functions, compliance, and business intelligence.
- To create insights to support and enable effective business analysis and decision making.



Why Use A Data Warehouse?

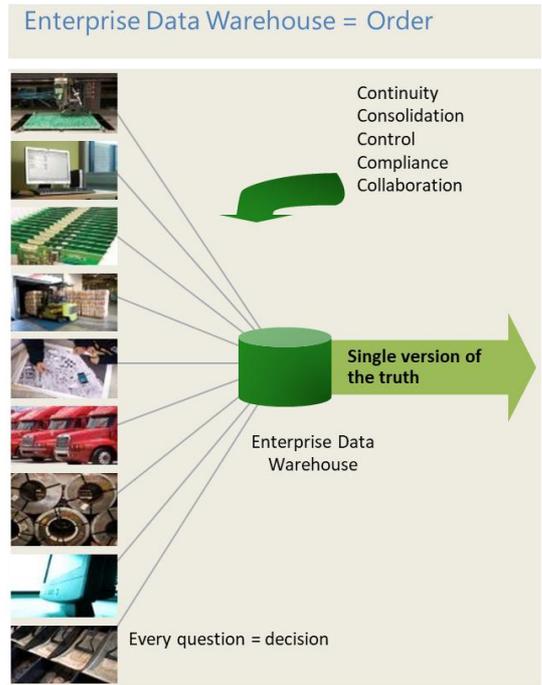
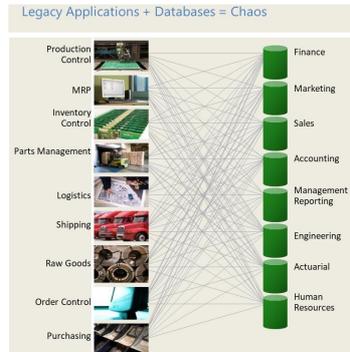
Legacy Applications + Databases = Chaos



Why Use A Data Warehouse?

Purposes of a Data Warehouse:

- 1) Standardize data from multiple sources;
- 2) Save time building reports & analyses;
- 3) Report & analyze in ways you could not do before;
- 4) Prevent adverse impact on Operational systems;
- 5) Enable comparison without the underlying data changing;
- 6)



Classic Characteristics of a DW (Bill Inmon)

"a subject oriented, integrated, time variant and non-volatile collection of summary and detailed historical data used to support the strategic decision-making processes for the corporation"

Subject-Oriented:

A data warehouse can be used to analyze a particular subject area. For example, "sales" can be a particular subject.

Integrated:

A data warehouse integrates data from multiple data sources. For example, source A and source B may have different ways of identifying a product, but in a data warehouse, there will be only a single way of identifying a product

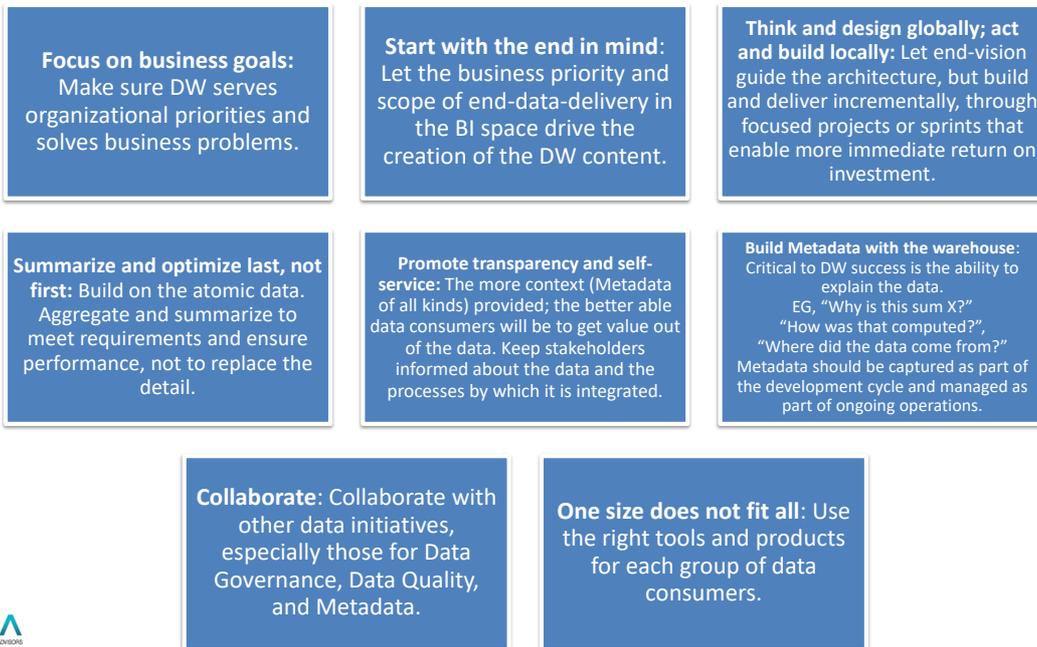
Time-Variant:

Historical data is kept in a data warehouse. For example, one can retrieve data from 3 months, 6 months, 12 months, or even older data from a data warehouse. This contrasts with a transactions system, where often only the most recent data is kept. For example, a transaction system may hold the most recent address of a customer, where a data warehouse can hold all addresses associated with a customer

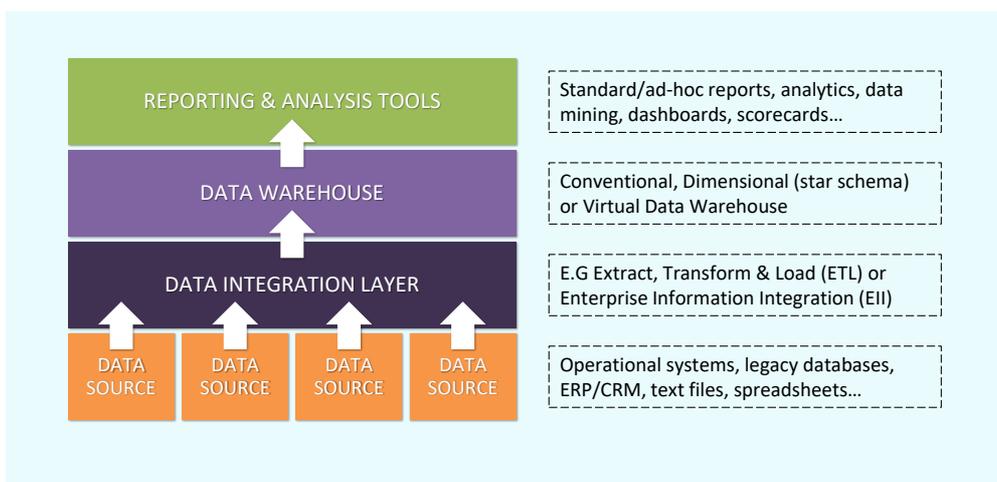
Non-volatile:

Once data is in the data warehouse, it will not change. So, historical data in a data warehouse should never be altered.

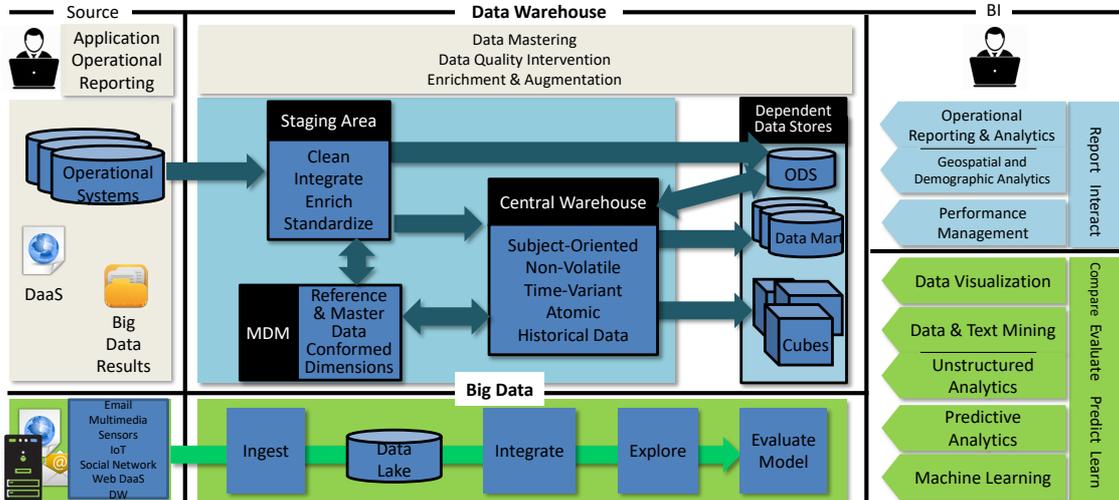
Data Warehouse implementation guiding principles



Simplified Business Intelligence Stack



Conceptual DW/BI and Big Data Architecture



	 Data Lake	 Data Warehouse	 Data Mart
Data Scope	Broad, Raw Data	General, Cleaned	Focused
Prior Processing	None – Light	Moderate - High	Very High
Analysis Limitations	Only limited by input sources	Limited by Data cleaning / transformation choices	Limited to the Mart topic focus
Ease of Navigation	Difficult	Moderate	Easy



Types of Data Warehousing

CLASSIC

for Strategic and Tactical BI

- Data is loaded to the DW by a batch process; typically nightly
- All data in the DW is non-volatile
- Allows what-if analyses & re-testing of hypotheses against a stable data set

ACTIVE

for Strategic, Tactical and Operational BI

- > Data, including volatile data, is loaded more frequently
- > Drivers: lower latency and the need for more real-time or near real-time data
 - » E.g. Real-time balances at a cash point (ATM)
- > Isolation of change
 - » Isolate volatile data from historical non-volatile data (e.g. put volatile data in an ODS)
 - » Use union queries or Data Virtualisation to combine volatile and non-volatile data.
- > Alternatives to batch ETL
 - » Data Virtualisation
 - » Message queues
 - » Trickle-feeds / Pipelining
 - » Event-based updates (Change Data Capture)



P / 182

Big Data

Early definition: 3 V's

Volume:

The amount of data. Big data often has thousands of entities or elements in billions of records.

Velocity:

The speed at which data is captured, generated, or shared. Big data is often generated and can also be distributed and even analyzed in real time.

Variety:

The forms in which data is captured or delivered.

Expanded definition: 6/7 V's

Viscosity:

How difficult the data is to use or integrate.

Volatility:

How often data changes occur and therefore how long the data is useful.

Veracity:

How trustworthy the data is.

Variability:

Big Data requires storage of multiple formats; data structure is often inconsistent within or across data sets.



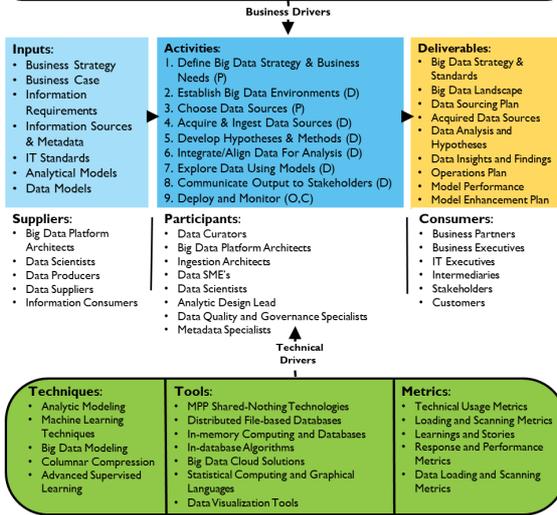
P / 183

Big Data and Data Science (DMBoK 2 revised)

Definition: The handling of large amounts of data (Big Data) and paradigm and statistical analytics (Data Science) of many different types of data to find answers and insights.

Goals:

1. Discover relationships between data and the business.
2. Discover and analyze new factors that might affect the business.
3. Package communications of model outputs for stakeholders and decision-makers.
4. Integrate existing organizational practices with best practices in data management, big data, and data science.



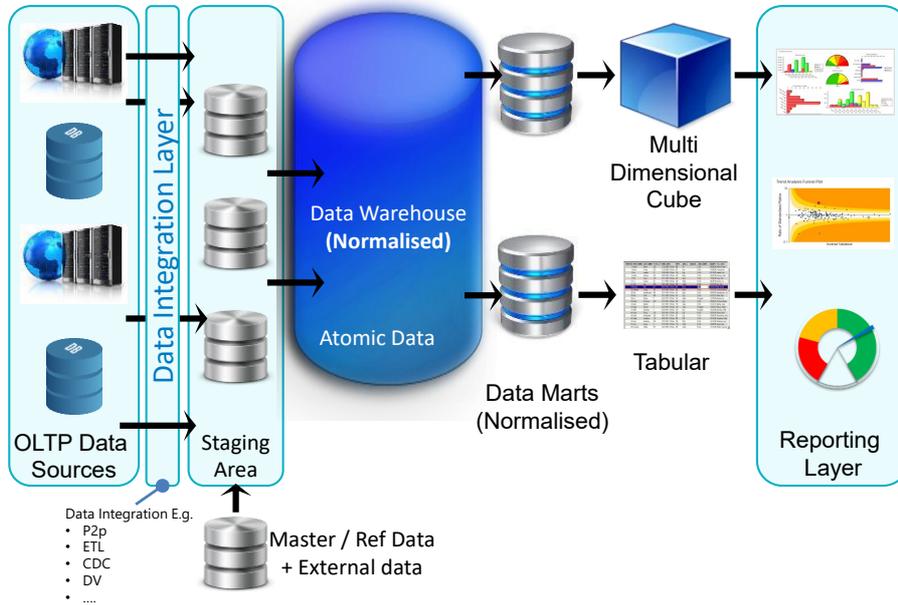
What is Data Warehousing? (DMBoK)

Data Warehousing is the term used to describe the processes that maintain the data contained within a data warehouse, namely:

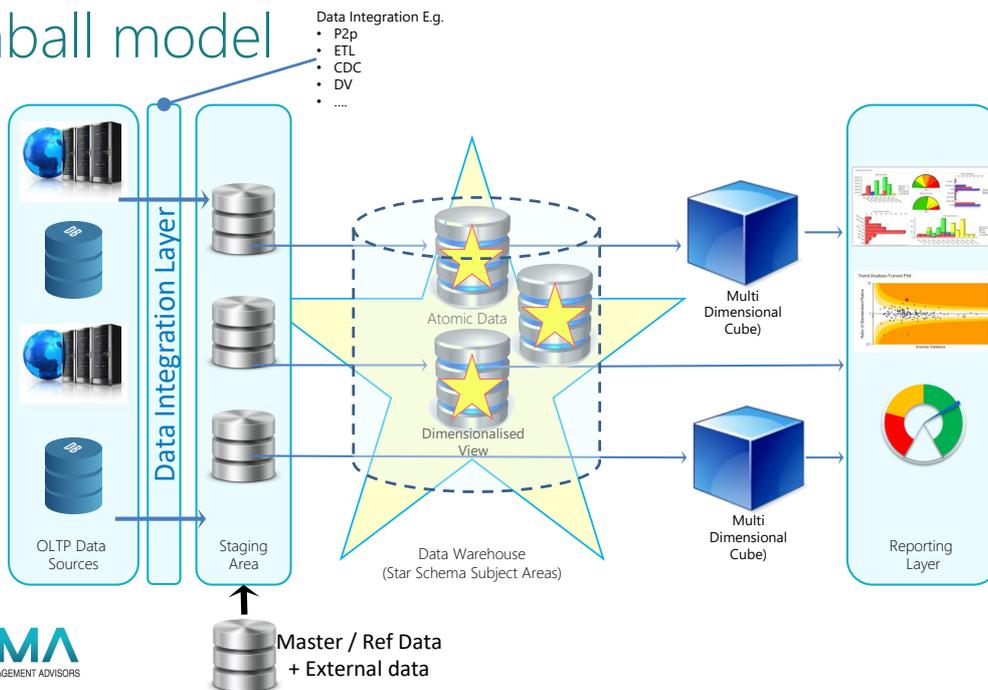
- Extract processes
- Cleansing processes
- Transformation processes
- Load processes
- Associated Control processes
- The use of Meta-data



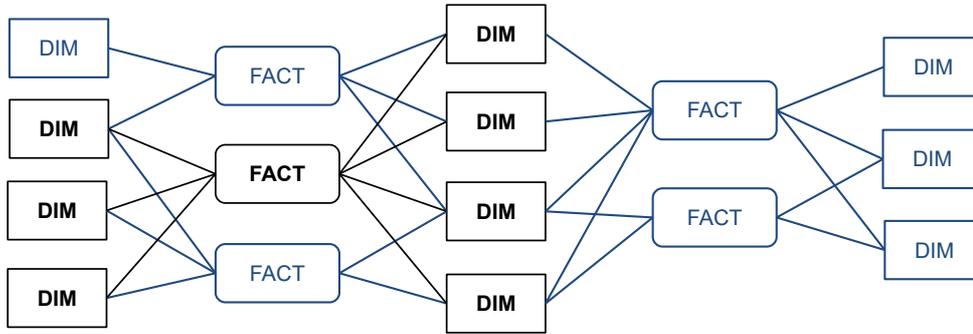
Inmon Data Warehouse Model



Kimball model



How do Dimensional Models fit into the Data Warehouse?



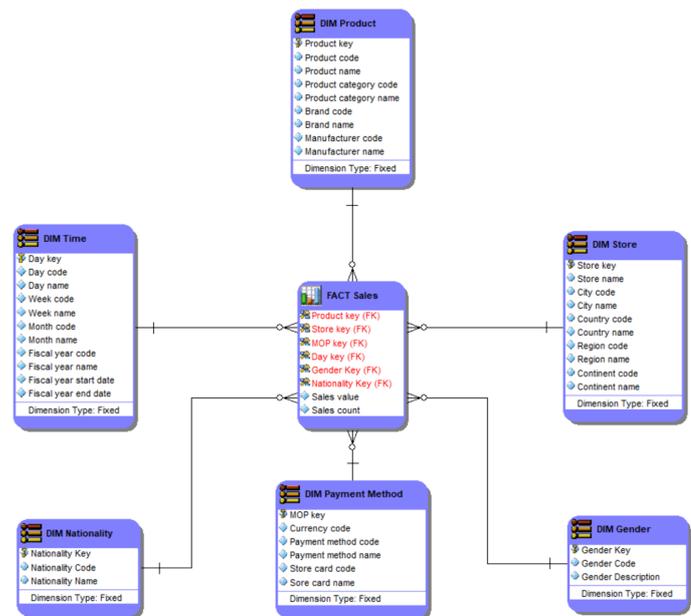
A Dimensional Model

Dimension tables

- Examples: Location, Product, Time, Promotion, Organisation etc.
- Records in the dimension tables correspond to nouns.
- The data in the dimension tables changes slowly – the number of new records created each day is typically low.

Fact tables

- Contains measures (e.g. Sales Value GBP) and dimension columns
- Records in the fact tables correspond to events, transactions, or measurements.
- The number of new records created each day is typically high.



Dimensions & Hierarchies

Hierarchies for the dimensions are stored in the dimensional table itself so there is no need for the individual hierarchical lookup tables be shown in the model.

Records in dimension tables correspond to nouns, the tables are "short" – 10s to 1,000s of records

DIM Product

- Product key
- Product code
- Product name
- Product category code
- Product category name
- Brand code
- Brand name
- Manufacturer code
- Manufacturer name

Dimension Type: Fixed

Rich set of attributes, tables are "wide" – many columns & the data changes slowly

Denormalised so no need to join to further lookup tables. This means there is redundancy

DIM Store

- Store key
- Store name
- City code
- City name
- Country code
- Country name
- Region code
- Region name
- Continent code
- Continent name

Dimension Type: Fixed



Fact Tables

Records in fact tables correspond to **events, transactions, or measurements.**

Data is added regularly; the tables are "long" – often millions of records

FACT Sales

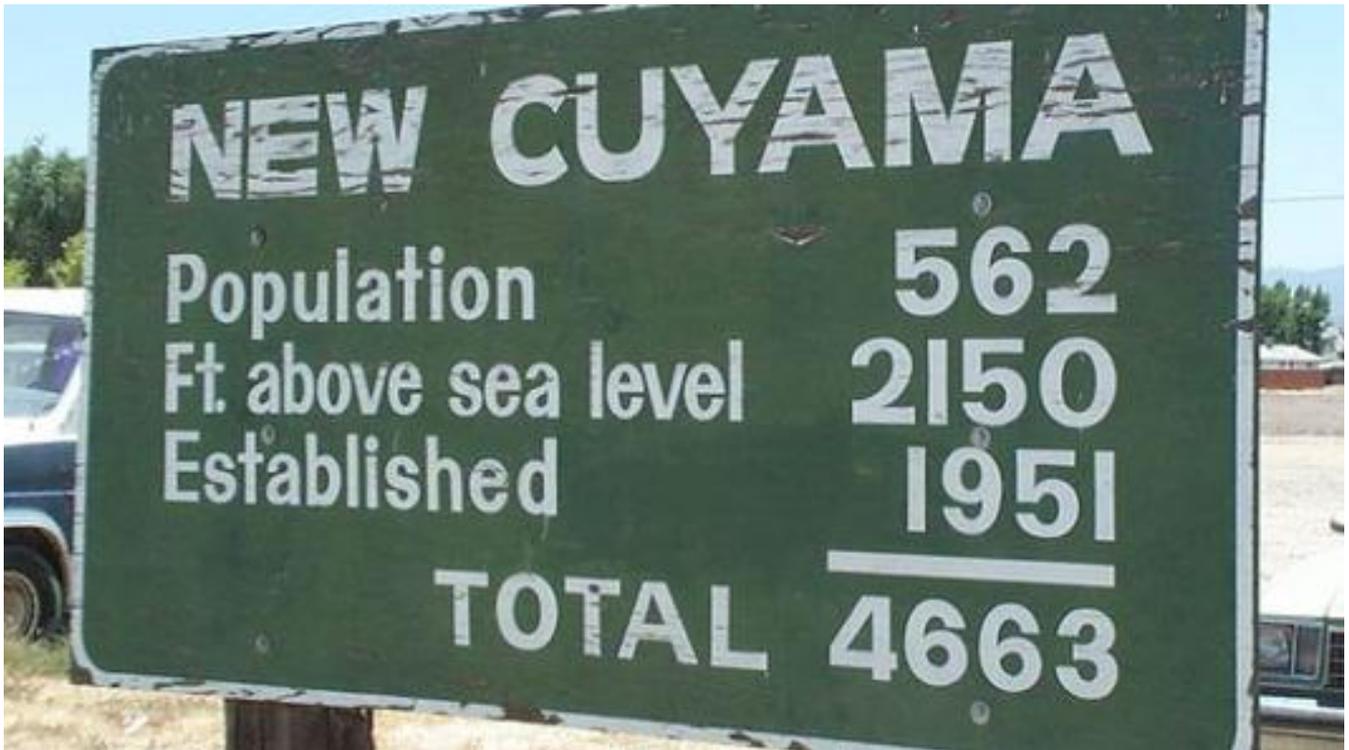
- Product key (FK)
- Store key (FK)
- MOP key (FK)
- Day key (FK)
- Gender Key (FK)
- Nationality Key (FK)
- Sales value
- Sales count

Rich set of attributes; the tables are "narrow" – minimal number of columns

Low redundancy

The most useful measures are "additive"





How can we handle slowly changing dimensions?



There are standard techniques for handling slowly changing dimensions.

- Type 1 (overwrite)
- Type 2 (add new row)
- Type 3 (add new attribute)
- Type 4 (add history table)
- Type 6 (hybrid)
- Others – see the internet!

We may need to employ different techniques for different fields.

Types of Business Intelligence Tool

Query and Reporting Tools	<ul style="list-style-type: none"> • Enable business users to create ad hoc queries and reports. Also used to create standard reports. • Examples: Business Objects, PowerBI, Cognos, Crystal Reports
OLAP Tools	<ul style="list-style-type: none"> • Support interactive multi-dimensional analysis • Example: Hyperion, MS Analysis Services with Excel
Analytic Applications	<ul style="list-style-type: none"> • “Industry Standard Data Marts in a Box” • Include logic and processes to extract data from well-known source systems (e.g. vendor ERP systems) • Also include a data model for the data mart and pre-built reports and dashboards.
Management Dashboards and Scorecards	<ul style="list-style-type: none"> • Dashboard = dynamic presentation of operational information • Scorecard = static representation of progress towards longer term goals
Performance Management Tools	<ul style="list-style-type: none"> • Tools to monitor performance (typically of the organisation) against goals. Such tools use data from source systems but also capture data directly, e.g., about goals and activities. • Applications include budgeting, planning and financial consolidation
Predictive Analytics and Data Mining	<ul style="list-style-type: none"> • Data Mining - tools that allow users to discover patterns, often using various algorithms, by exploring data interactively, without having a specific question to start with. • Predictive analytics – tools that predict the future and support “what if” analysis.
Advanced Visualisation and Discovery Tools	<ul style="list-style-type: none"> • Use an in-memory architecture to allow users to interact with the data in a highly visual, interactive way. • Example: PowerBI, Tableaux, QlikView

Summary

- A Data Warehouse is a technology solution that addresses the problem of how to support **business intelligence** needs.
- Data Warehouses contain data from source systems together with externally provided data that is combined & stored in a manner optimised for **query intensive use**
- Data Warehouses must contain **integrated** data in order to be more than just a warehouse of data.
- Dimensional Data Modelling is used to make data **easy to use** and to support **fast queries**.
- Data Warehousing processes must be transparent (open) and visible in order for business users to **understand** and have **confidence** in data.
- **Many** types of **Business Intelligence tools** exist

Ref	Question	A	B	C	D	E
DW1	Which of the following uses for a Data Mining tool is not optimal?	Identification of data quality issues with your SAP Financial system	Fraud Detection	Customer Segmentation and Scoring	Predictive Analysis	Identifying potential loan defaulter's
DW2	Which of the following is not a good example of BI?	Strategic Analytics for Business Decisions	Decision Support Systems	Supporting Risk Management Decision Reporting	Statutory reporting to a Regulatory Body	Identifying top quartile customers
DW3	According to Henry Morris of IDC, Analytic Applications provide business with a pre-built solution to optimize a functional area or industry segment	TRUE	FALSE			
DW4	When performing an evaluation of analytic applications, which of the following questions is least relevant to identify the level of effort needed?	How much of the tool infrastructure meets our organisational infrastructure	The Standard source systems for which ETL is supplied	No. of source systems we need to integrate into the tool	How much do the canned processes in the tool match our business	Annual costs such as license, maintenance, etc
DW5	You need to discover possible relationships or to show data patterns in an exploratory fashion when you do not necessarily have a specific question to ask. What kind of data tool would you use to identify patterns of data using various algorithms?	ETL Jobs	Data Quality Profiler	Meta-Data Data Lineage View	Data Mining	Data Visualisation Application
DW6	"Slice", "Dice", "Roll-up" and "Pivot" are terms used in what kind of data processing?	OLAP	OLTP	ODS	EDI	EIEIO
DW7	A comparatively new architectural approach is where volatile data is provisioned in a data warehouse structure to provide transactional systems with a combination of historical and near real time data to meet customer needs. This is a definition of:	Operational Data Store	Behavioural Decision Support Systems	Active Data Warehousing	On Line Transactional Processing System	On Line Analytical Processing Cube

AFTER QUIZ 5

1. General (6)
2. Data Quality (9)
3. Data Storage & Operations (3)
4. Master & Reference Data (9)
5. DW / BI + Big Data (7)

Maximum possible score = 34

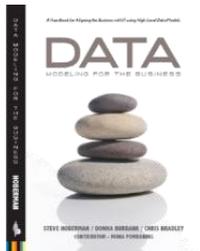
60% (CDMP Associate) = 21

70% (CDMP Practitioner) = 24

80% (CDMP Master) = 28

Data Modelling & Design

Data Management Process	2%
Big Data	2%
Data Architecture & Lifecycle Management	6%
Document, Records & Content Management	4-5%
Data Ethics	2%
Data Governance	11%
Data Integration & Interoperability	6-7%
Master & Reference Data Management	11%
Data Modelling & Design	11%
Data Quality	11%
Data Security	6%
Data Storage & Operations	4-5%
Data Warehousing & Business Intelligence	11%
Metadata Management	11%

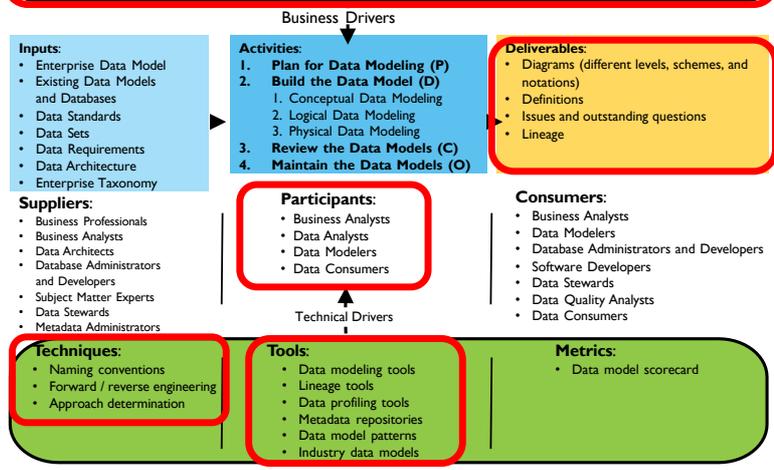


Data Modeling and Design DMBOK2 Revised

Definition: Data modeling is the process of discovering, analyzing, and scoping data requirements, and then representing and communicating these data requirements in a precise form called the data model. This process is iterative and involves conceptual, logical, and physical models

- Goals:**
- To confirm and document an understanding of different perspectives, which leads to applications that more closely align with current and future business requirements.
 - To understand how data fits together and creates a foundation to successfully complete broad-scope initiatives such as master data management and data governance.
 - To lower support costs and reduce the costs of building new applications thanks to increased reusability opportunities.

“Data Modelling For The Business
 A Handbook for aligning the business with IT using high-level data models”
 Technics Publishing
 ISBN 978-0-9771400-7-7;
<http://www.amazon.com/Data-Modeling-Business-Handbook-High-Level>



(P) Planning,
 (C) Control,
 (D) Development,
 (O) Operations



The deliverables of the data modeling process include:

- **Diagram:** A data model contains one or more diagrams. The diagram is the visual that captures the requirements in a precise form. It depicts a level of detail (e.g., conceptual, logical, or physical), a scheme (relational, dimensional, object-oriented, fact-based, time-based, or NoSQL), and a notation within that scheme (e.g., information engineering, unified modeling language, object-role modeling).
- **Definitions:** Definitions for entities, attributes, and relationships are essential to maintaining the precision on a data model.
- **Issues and outstanding questions:** Frequently the data modeling process raises issues and questions that may not be addressed during the data modeling phase. In addition, often the people or groups responsible for resolving these issues or answering these questions reside outside of the group building the data model. Therefore, often a document is delivered that contains the current set of issues and outstanding questions. An example of an outstanding issue for the student model might be, “If a **Student** leaves and then returns, are they assigned a different **Student Number** or do they keep their original **Student Number**?”
- **Lineage:** For physical and sometimes logical data models, it is important to know the data lineage, that is, where the data comes from. Often lineage takes the form of a source/target mapping, where one can capture the source system attributes and how they populate the target system attributes. Lineage can also trace the data modeling components from conceptual to logical to physical within the same modeling effort. There are two reasons why lineage is important to capture during the data modeling. First, the data modeler will obtain a very strong understanding of the data requirements and therefore is in the best position to determine the source attributes. Second, determining the source attributes can be an effective tool to validate the accuracy of the model and the mapping (i.e., a reality check).

Page 149 / 152 in printed version
Page 156-157 of DM-BOK pdf
Section 2.1 (Activities / Plan for
Data Modelling)

P / 201

What is a Data Model?

A model is a representation of something in our environment making use of standard symbols to enable improved understanding of the concept

A data model describes the specification, definition and rules for data in a business area

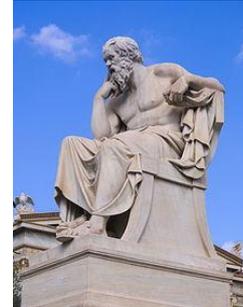
A data model is a diagram (with additional supporting metadata) that uses text and symbols to represent data to give the reader a better understanding of the data

A data model describes the inherent logical structure of the data within a given domain and, by implication, the underlying structure of that domain itself



Clarity & Definition is vital

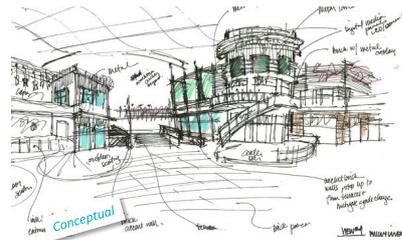
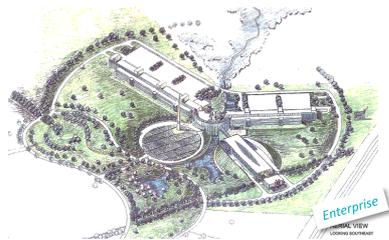
“The beginning of wisdom is the definition of terms.”



Socrates: 470 – 399 BC

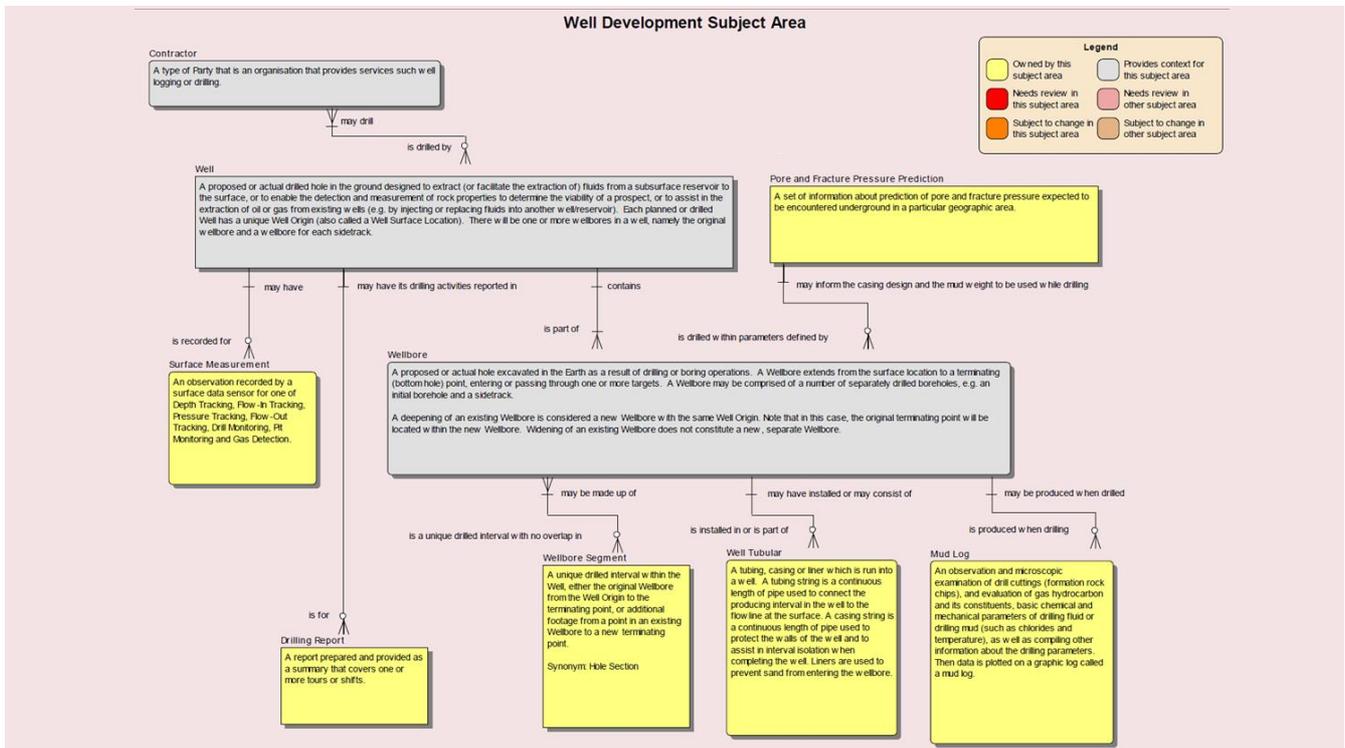
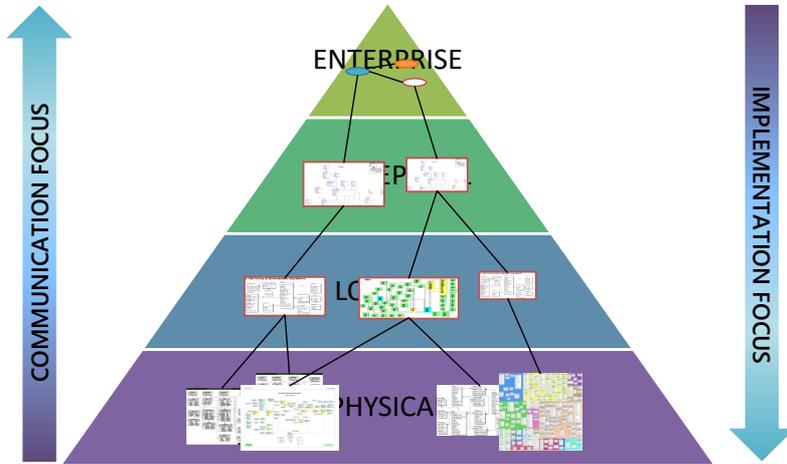


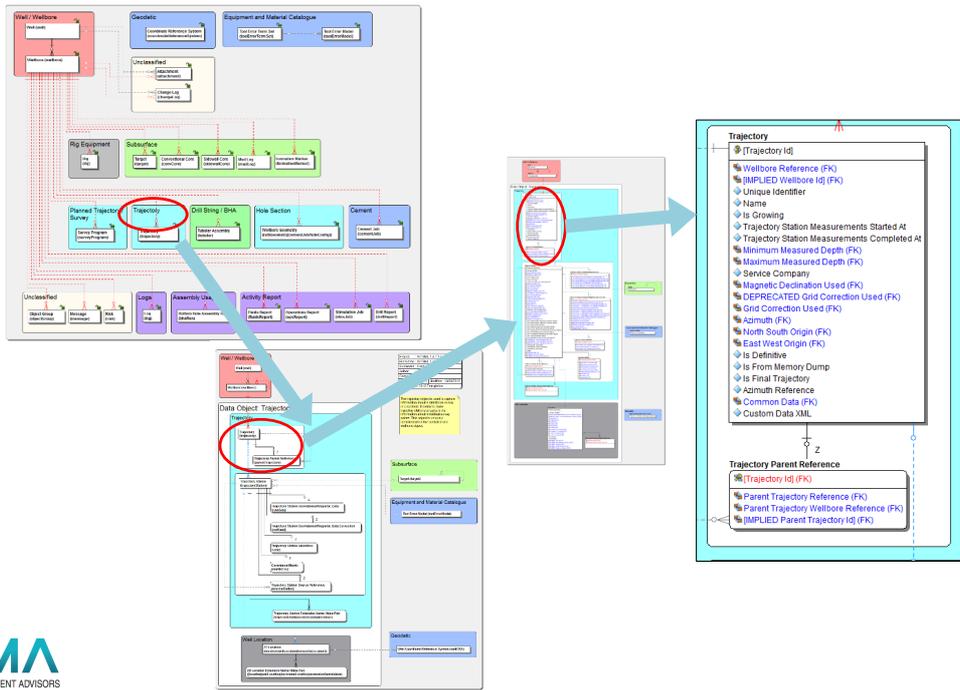
We All Use Models



Data Model Levels

ANSI Standards Planning and Requirements Committee (SPARC) call these levels: Conceptual, External, and Internal



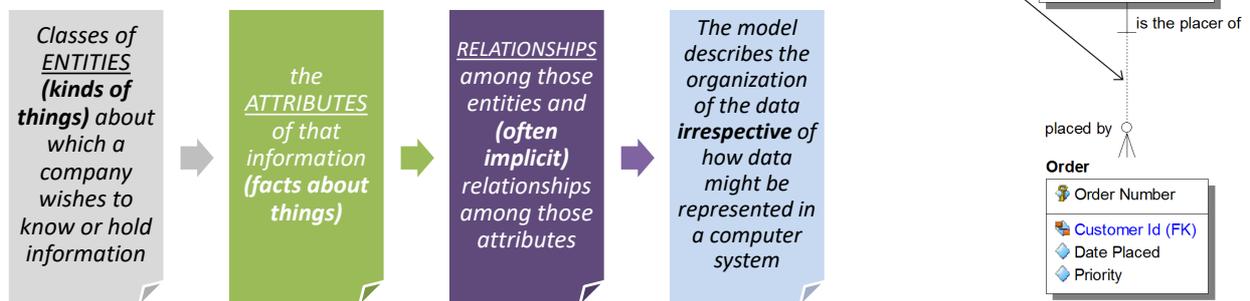




The main purpose of a data model is not to design a database, but to describe a business

(C. Bradley CDMP Fellow 2010)

A Data Model represents



A Data Model Represents (Entity Categories)

Classes of **entities** (kinds of things) about which a company wishes to know or hold information

- WHO** *Person, Employee, Vendor, Customer, Department, Organisation, ...*
- WHAT** *Product, Service, Raw Material, Training Course, Flight, Room, ...*
- WHEN** *Time, Day, Date, Calendar, Reporting Period, Fiscal Period, ...*
- WHERE** *Geographic location, Delivery address, Storage Depot, Airport, ...*
- WHY** *Order, Complaint, Inquiry, Transaction, ...*
- HOW** *Invoice, Policy, Contract, Agreement, Document, Account, ...*

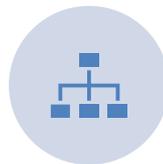


Types of data modelled (DMBoK)



Category information:

Data used to classify and assign types to things. For example, customers classified by market categories or business sectors; products classified by color, model, size, etc. AKA Reference Data.



Resource information:

Basic profiles of key resources required to conduct operational processes such as Product, Customer, Supplier, Facility, Organization, and Account. AKA Master & Reference Data.



Business event information:

Data created while operational processes are in progress. Examples include Customer Orders, Supplier Invoices, Cash Withdrawal, and Business Meetings. AKA transactional business data.



Detail transaction information:

Often produced through point-of sale systems (either in stores or online). Also produced through social media systems, other Internet interactions (clickstream, etc.), and by sensors in IOT / machines (GPS, RFID, Wi-Fi, etc.). Often it is aggregated, used to derive other data, and analyzed for trends. This type of data (large volume and/or rapidly changing) is often referred to as Big Data.



What Is A Business Data Model?

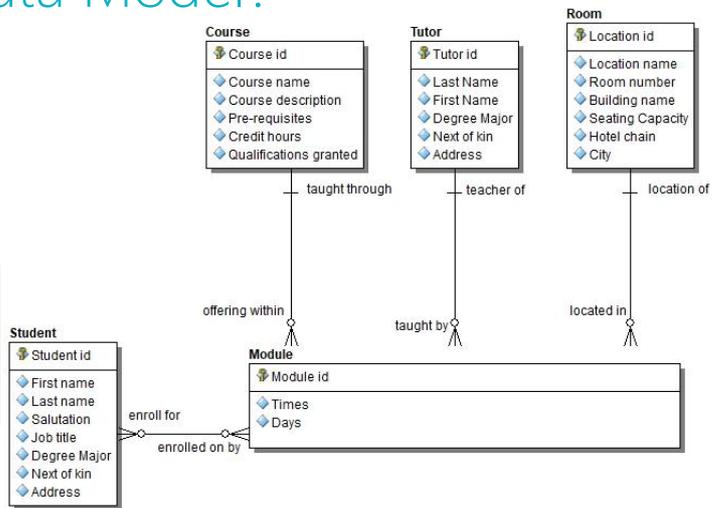
- A description of a Business (or an area of the Business) in terms of the **things** it needs to know about.
- The Data **things** are “entities” and the “**facts about things**” are attributes & relationships.
- It’s a representation of the “real world”, not a technical implementation of it
- *Should* be able to be understood by Business users

Definition:

A Student is any person who has been admitted to a course, has paid, and has enrolled in one or more modules within a course. Tutors and other staff members may also be Students

Business Assertions

- A Student enrolls for zero one or more modules
- A Course can be taught through zero one or more Modules
- A Room can be the location of zero one or more modules
- A Tutor can be the teacher of zero one or more modules



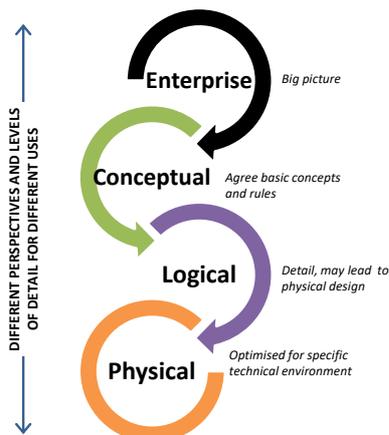
The Other Way?

- A Module is enrolled in by many students
- A Module is an offering within one course
- A Module is located in one room
- A Module is taught by one tutor

Really?



Enterprise vs. Conceptual vs. Logical

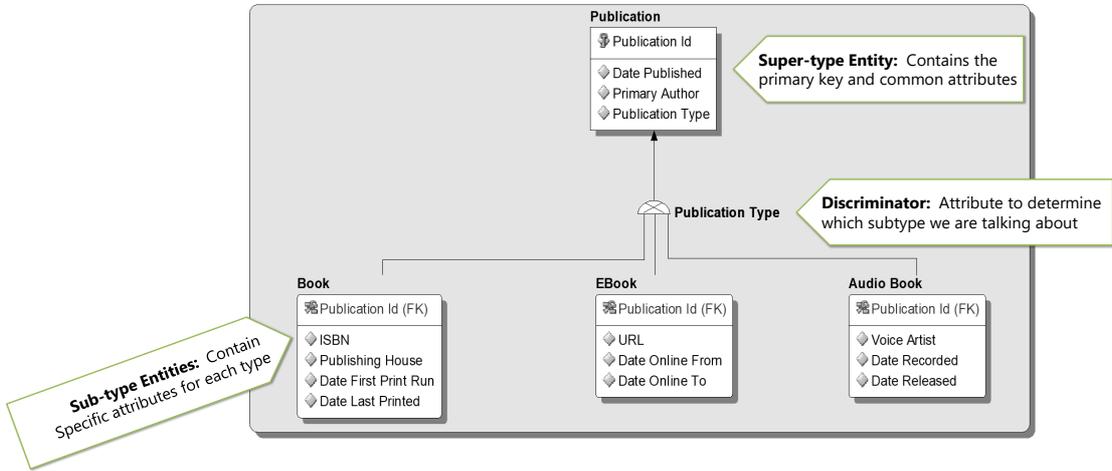


- Common understanding before progressing too far into detail
- Used to communicate with the Business
- Overview: main entities, super types, attributes, and relationships
- Lots of Many to Many & multi meaning relationships
- Relationships frequently show multiplicity of meaning
- May be denormalised
- Non-atomic & multi-valued attributes allowed; no keys
- Should fit on one page
- 20% of the modelling effort

- Detailed: ~ 5x Entities vs Conceptual model
- Detailed: Frequently pre-cursor to 1st cut physical (database) design
- Detailed: Key input to requirements specification
- M:M relationships resolved: Intersection entities mostly have meaning
- Supertypes AND Subtypes included
- Relationship optionality added
- Primary, foreign, alternate keys included
- Reference entities included
- Fully normalized – no multi-valued, redundant, non-atomic attributes
- May be partitioned (sub-models)
- 80% of the modelling effort

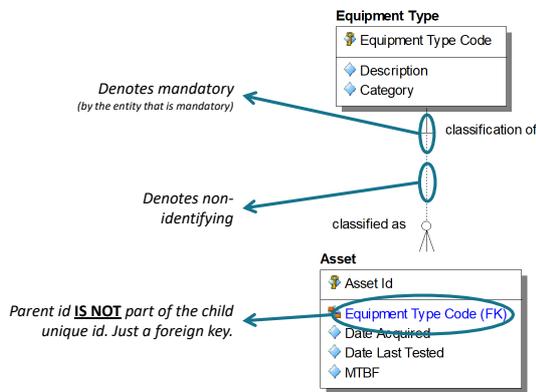


Entity Subtypes



Relationship Types

Non-Identifying



Eq Type	Eq Type Code	Asset Id
Laptop	L	1
Laptop	L	2
Laptop	L	3
Projector	P	4
Projector	P	5
Desktop	D	6
Desktop	D	7
Desktop	D	8
Desktop	D	9

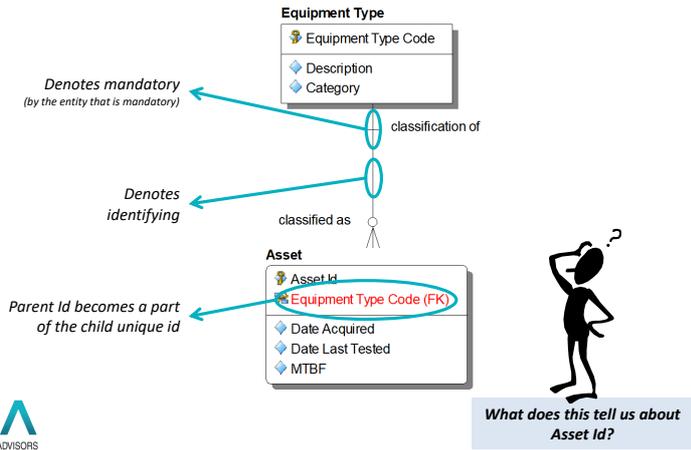


What does this tell us about Asset Id?



Relationship Types

Identifying

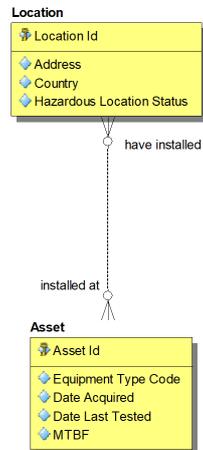
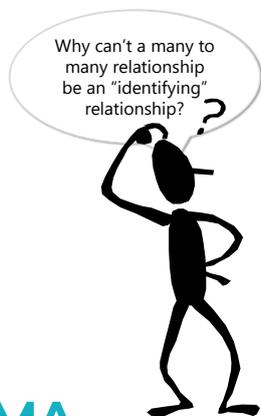


Eq Type	Eq Type Code	Asset Id
Laptop	L	1
Laptop	L	2
Laptop	L	3
Projector	P	1
Projector	P	2
Desktop	D	1
Desktop	D	2
Desktop	D	3
Desktop	D	4



Relationship Types

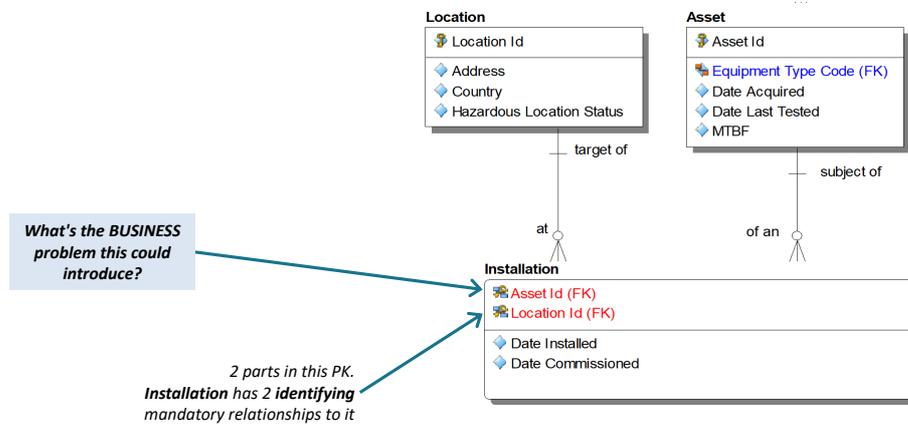
Many to Many (AKA nonspecific)



Chris's law
99% of M:M relationships represent a real business concept that is the intersection entity

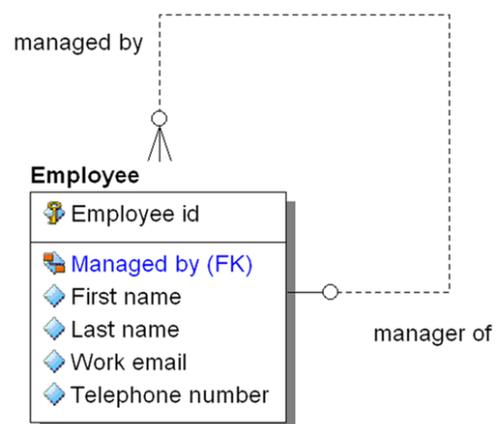


Relationship Types

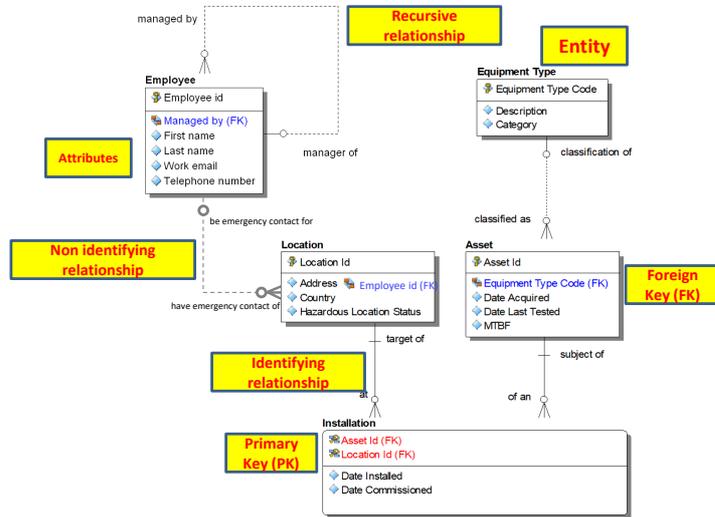


Recursive Relationships

- A recursive relationship occurs when there is a relationship **between an entity and itself**.
- E.g. A one-to-many recursive relationship occurs when an **employee** is the **manager** of other **employees**.
- The **employee** entity is related to itself, and there is a one-to-many relationship between one **employee** (the manager) and many other **employees** (the people who report to the manager).
- **But we cannot have duplicated attribute names in an entity; hence the FK is given a role name (e.g. "managed by").**

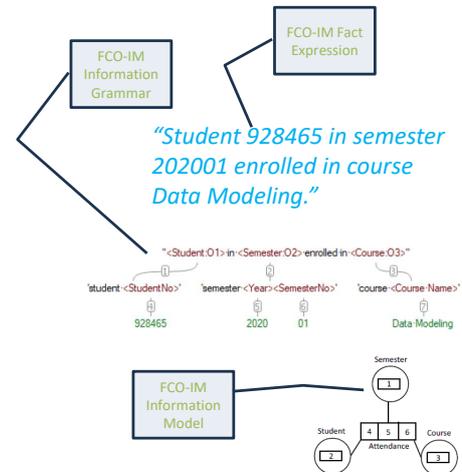


Logical Data Model Components

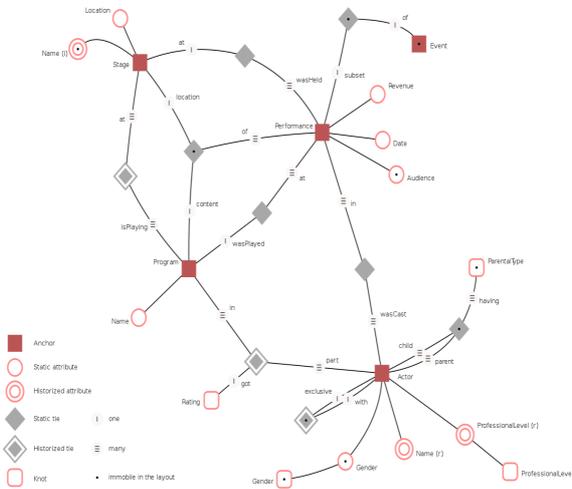


Overview of Data Modelling Notations

- **Entity-Relationship (ER) Notation:** For example, Information Engineering (IE), uses a set of symbols for representing entities, their attributes, and relationships
- **Barker Notation:** Developed by Richard Barker, used in the Oracle Designer CASE tool and is popular in Europe
- **IDEF1X Notation:** Used in the integration definition for information modelling, a popular Federal standard for data modelling
- **Fully Communication Oriented Modeling (FCO-IM):** The conceptual information models are made using natural language with concrete examples. Using so called elementary fact expressions, Information grammar & Information model; FCO-IM provides extensive support for different ways of communicating similar facts, synonyms, homonyms, generalized object types, and other conceptual modeling aspects.
- **UML Class Diagram Notation:** Used in the Unified Modelling Language for representing object-oriented software systems
- **Data vault modelling:** A database modelling technique that provides historical storage of data from multiple operational systems.



Data Vault & Anchor Models



Well suited to handle information that changes in structure and context over time.

Hubs, Satellites, and Links are the 3 entity types to focus its DV design around functional areas of business.

Data Vault Modeling is a hybrid approach, encompassing the best of breed between 3NF and Dimensional (star schema)

Anchors are shown as squares



Keys



Key Vocabularies

Terms used conceptually

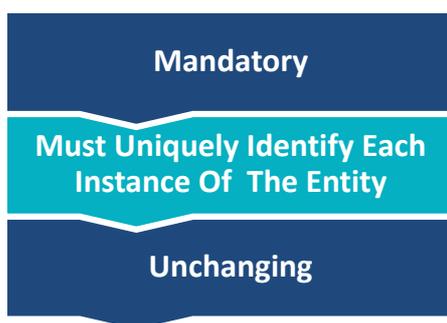
- Primary key
- Alternate key
- Composite key
- Super key
- Candidate key
- Surrogate key
- ...

Terms used physically

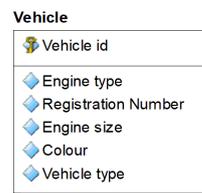
- Primary key
- Clustered key
- Encryption key
- Partitioning key
- Index
- Identity
- Sequence
- ...



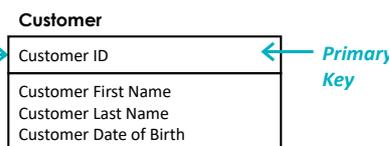
Primary Key



Can I use the Vehicle's Registration Number as the PK?



Attribute(s) that make up a PK are represented in modelling tools separately from the rest of the attributes by a line



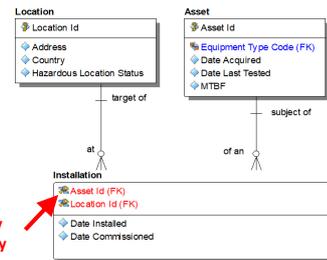
Composite vs Compound (DMBoK2_R)

A **composite key** is a candidate key that consists of two or more attributes that together uniquely identify an entity.

Student	Date	Period	Present
J. Smith	01/06/09	1	Y
B. Rogers	01/06/09	1	Y
A. Black	01/06/09	2	N
H. Rose	02/06/09	1	Y
M. Wright	02/06/09	2	N

These three fields are combined to make a composite primary key

A **compound key** is a composite key for which each attribute that makes up the key is a foreign key in its own right

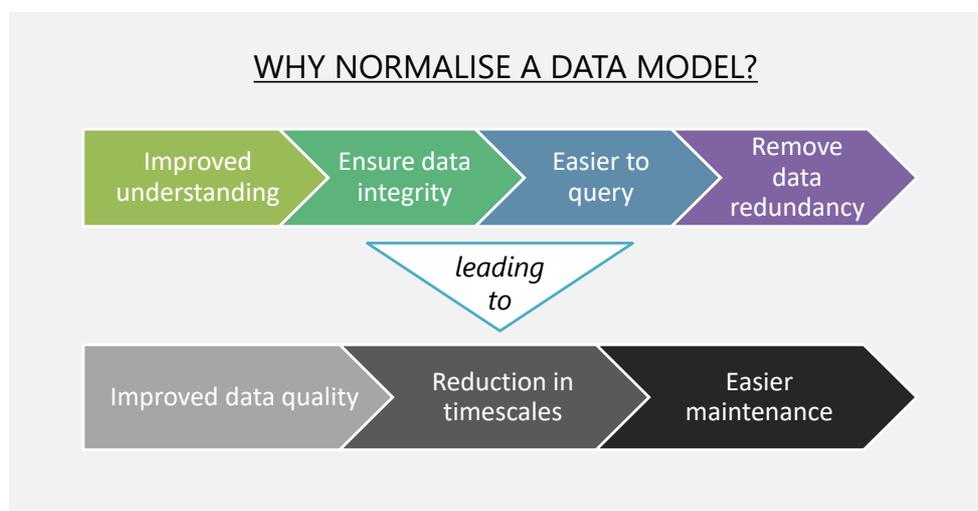


These 2 foreign keys make up the primary key of the entity and are a Compound primary key

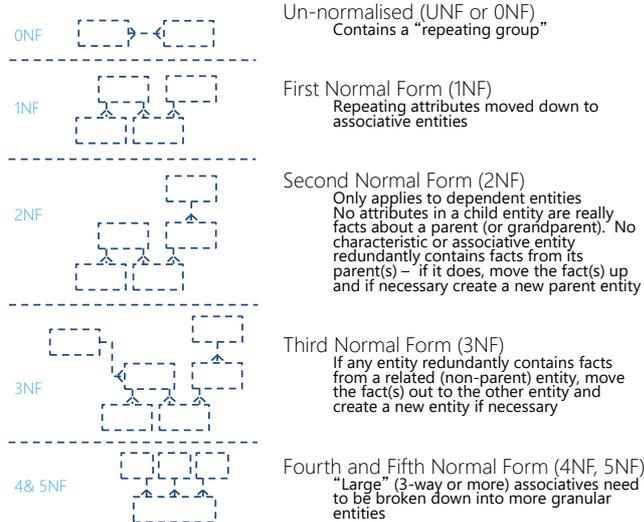
A compound key is similar to a composite key in that two or more attributes are needed to create a unique value. However, a compound key is created when two or more primary keys from different entities are present as foreign keys within an entity. The foreign keys are used together to uniquely identify the entity.



Normalisation



Summary

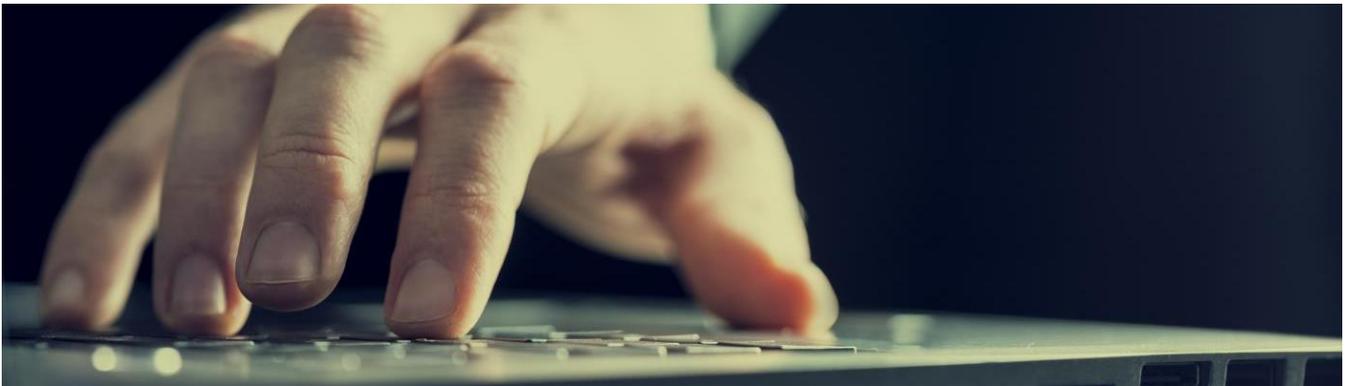


Definition

1NF Every non-key attribute in an entity must depend on it's primary key

2NF Each entity must have the fewest possible correct primary key attributes

3NF Each non key element must be directly dependent upon the primary key and not upon any other non-key attributes



Physical database design best-practice

- Use normalised design for relational databases supporting OLTP apps.
- Use views, functions and stored procedures to create non-normalised, application-specific, object-friendly, conceptual (virtual) views of data.
- Use standard naming conventions.
- Enforce data security and integrity at the database level, not in the application.
- Keep database processing on the database server as much as possible.
- Grant permissions on database objects only to application groups or roles, not to individuals.
- Do not permit any direct, ad-hoc updating of the database.



Class, Prime, Modifier, Qualifier Words

The following word classification types are used by various data modelling tools and are defined below with examples.

Defining Word Classification Types

Prime Word:

The prime word identifies the object or element being defined. Typically, these objects represent a person, place, thing, or event about which an organization wishes to maintain information. Prime words may act as primary search identifiers when querying a database system and provide a basic list of keywords for developing a general-to-specific classification scheme based on business usages. **CUSTOMER** in **Customer Address** is an example of a **prime word**.

Modifier:

A modifier gives additional information about the class word or prime word. Modifiers may be adjectives or nouns. **DELIVERY** in **Customer Delivery Address** is an example of a **modifier**. Other modifier examples: **ANNUAL**, **QUARTERLY**, **MOST**, and **LEAST**.

Class Word:

A class word is the most important noun in a data element name. Class words identify the use or purpose of a data element. Class words designate the type of information maintained about the object (prime word) of the data element name. **ADDRESS** in **Customer Address** is an example of a **class word**.

Qualifier:

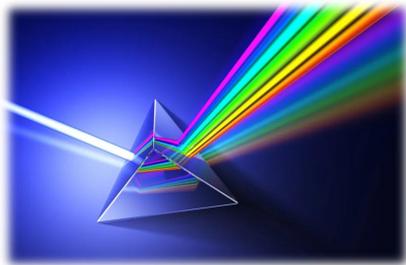
A qualifier is a special kind of modifier that is used with a class word to further describes a characteristic of the class word within a domain of values, or to specify a type of information that can be attached to an object.

Examples: FEET, METERS, SECONDS, and WEEKS.



P / 234

Physical database design best practices



Performance and Ease of Use

Ensure quick and easy access to data

Reusability

Multiple applications can use the data

Integrity

The data should have valid business meaning and value

Security

Data should only be available to authorised users

Maintainability

Ensure cost of maintenance does not exceed its value to the organisation



P / 235

Transforming from a logical to physical data model



- **Denormalisation**
Selectively and justifiably violating normalisation rules to reduce retrieval time, potentially at the expense of additional space, insert/update time and reduced data quality.
- **Surrogate keys**
Substitute keys not visible to the business.
- **Indexing**
Create additional index files to optimise specific types of queries.
- **Partitioning**
Break a table or file horizontally or vertically.
- **Views**
Virtual tables used to simplify queries, control data access and rename columns.
- **Dimensionality**
Creation of fact tables with associated dimension tables. Structured as star schemas and snowflake schemas for BI.

Partitioning



Horizontal partitioning

Horizontal partitioning is the partitioning of a table into a number of smaller tables on the basis of rows. For example, in an employee table, employees with a salary of less than £25, 000 will be partitioned into a different table.

Vertical partitioning

Vertical partitioning is dividing the table based on the different columns. For example, in a customer table, retrieving only the name and contact number of customers into a different table.

ACID Test For Transaction Processing

ATOMICITY

Atomicity requires that database modifications must follow an "all or nothing" rule. Each transaction is said to be atomic. If one part of the transaction fails, the entire transaction fails and the database state is left unchanged.

To be compliant with the 'A', a system must guarantee the atomicity in each and every situation, including power failures / errors / crashes.

This guarantees that 'an incomplete transaction' cannot exist.

CONSISTENCY

The consistency property ensures that any transaction the database performs will take it from one consistent state to another.

Consistency states that only consistent (valid according to all the rules defined) data will be written to the database.

Quite simply, whatever rows will be affected by the transaction will remain consistent with each and every rule that is applied to them (including but not only: constraints, cascades, triggers).

While this is extremely simple and clear, it's worth noting that this consistency requirement applies to everything changed by the transaction, without any limit (including triggers firing other triggers launching cascades that eventually fire other triggers etc.) at all.

ISOLATION

The requirement that no transaction should be able to interfere with another transaction at all. In other words, it should not be possible that two transactions affect the same rows run concurrently, as the outcome would be unpredictable, and the system thus made unreliable.

The ISOLATION property of ACID is often relaxed (i.e., partly respected) because of the speed decrease this type of concurrency management implies. In effect the only strict way to respect the isolation property is to use a serial model where no two transactions can occur on the same data at the same time and where the result is predictable (i.e. transaction B will happen after transaction A in every single possible case).

In reality, many alternatives are used due to speed concerns, but none of them guarantee the same reliability.

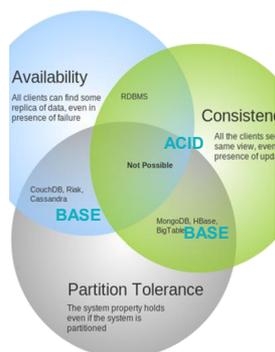
DURABILITY

Durability means that once a transaction has been committed, it will remain so. In other words, every committed transaction is protected against power loss/crash/errors and cannot be lost by the system and can thus be guaranteed to be completed.

In a relational database, for instance, once a group of SQL statements execute, the results need to be stored permanently. If the database crashes right after a group of SQL statements execute, it should be possible to restore the database state to the point after the last transaction committed.



BaSE



CAP computer science theorem quantifies the trade-offs.

[Eric Brewer's CAP theorem:](#)

If you want consistency, availability, and partition tolerance, you have to settle for two out of three.

For a distributed system, **partition tolerance** means the system will continue to work unless there is a total network failure. A few nodes can fail and the system keeps going.

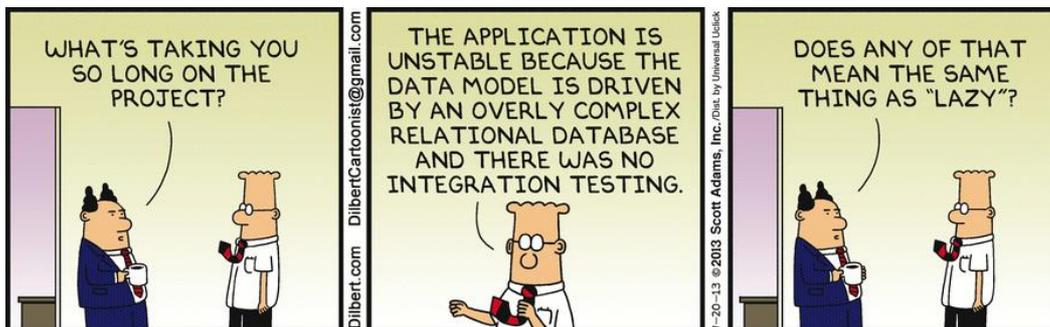
BASE is an alternative to **ACID**:

- **B**asic Availability
- **S**oft-state
- **E**ventual consistency

Rather than requiring consistency **after every transaction**, it is enough for the database to **eventually be in a consistent state**. (*Accounting systems do this all the time ... called "closing out the books."*)

It's OK to use stale data, and it's OK to give approximate answers.





Ref	Question	A	B	C	D	E
DD1	What are relationship labels?	The verb phrases describing the business rules in each direction between two entities.	The nullability setting on a foreign key.	A foreign key that has been role-named.	A relationship without cardinality.	A non-identifying relationship.
DD2	Which of these statements has the most meaningful relationship label?	An order line contains orders	An order is associated with order lines.	An order is related to order lines.	An order is connected with order lines.	An order is composed of order lines.
DD3	All of the following are TRUE statements on relationship types except:	A many-to-many relationship says that an instance of each entity may be associated with many instances of the other entity, and vice versa.	A one-to-one relationship says that a parent entity may have one and only one child entity.	A one-to-many relationship says that a parent entity may have one or more child entities.	A one-to-many relationship says that a child entity may have one or more parent entities.	A recursive relationship relates instances of an entity to other instances of the same entity.
DD4	In the BASE vs ACID model for Transaction Processing, E is best described which of these statements?	End to End data consistency	Eventual Availability of Data as described by the CAP theorem	Eventual Data Consistency	Extra Validation	Business Availability of Secure data Elements
DD5	A bank applies the business rule that each Customer may own one or many Accounts and each Account must be owned by one or many Customers. Which relationship type would be most appropriate?	many-to-many.	one-to-many.	many-to-one.	one-to-one.	recursive.
DD6	An employee may work for one other employee and may manage one or more employees. There is an indeterminate number of levels in this management hierarchy. What type of relationship would work best?	recursive.	one-to-one.	subtyping.	identifying.	non-identifying.
DD7	Which of the following business rules should NOT appear on a logical data model?	Each Person can work for zero to many Companies	Customer Last Name requires a non-unique index to improve retrieval performance.	Each Company must employ one or many Persons.	Each Order can contain one or many Order Lines.	Each Policy must belong to one Policy Owner.
DD8	All of the following are properties of a logical data model except:	contains relationship cardinality	technology-independent	contains attributes.	contains primary keys.	technology-dependent
DD9	According to the DMBOK, the deliverables for a data modelling process do not have to include:	the steps in the business process that use the data.	one or more diagrams.	definitions for entities, attributes and relationships.	issues and outstanding questions regarding data usage.	lineage describing where the data came from.

The deliverables of the data modeling process include:

- Diagram:** A data model contains one or more diagrams. The diagram is the visual that captures the requirements in a precise form. It depicts a level of detail (e.g., conceptual, logical, or physical), a scheme (relational, dimensional, object-oriented, fact-based, time-based, or NoSQL), and a notation within that scheme (e.g., information engineering, unified modeling language, object-role modeling).
- Definitions:** Definitions for entities, attributes, and relationships are essential to maintaining the precision on a data model.
- Issues and outstanding questions:** Frequently the data modeling process raises issues and questions that may not be addressed during the data modeling phase. In addition, often the people or groups responsible for resolving these issues or answering these questions reside outside of the group building the data model. Therefore, often a document is delivered that contains the current set of issues and outstanding questions. An example of an outstanding issue for the student model might be, "If a **Student** leaves and then returns, are they assigned a different **Student Number** or do they keep their original **Student Number**?"
- Lineage:** For physical and sometimes logical data models, it is important to know the data lineage, that is, where the data comes from. Often lineage takes the form of a source/target mapping, where one can capture the source system attributes and how they populate the target system attributes. Lineage can also trace the data modeling components from conceptual to logical to physical within the same modeling effort. There are two reasons why lineage is important to capture during the data modeling. First, the data modeler will obtain a very strong understanding of the data requirements and therefore is in the best position to determine the source attributes. Second, determining the source attributes can be an effective tool to validate the accuracy of the model and the mapping (i.e., a reality check).

Page 149 / 152 in printed version
Page 156-157 of DM-BOK pdf
Section 2.1 (Activities / Plan for Data Modelling)

AFTER QUIZ 6

1. General (6)
2. Data Quality (9)
3. Data Storage & Operations (3)
4. Master & Reference Data (9)
5. DW / BI + Big Data (7)
6. Data Modelling (9)

Maximum possible score = 43

60% (CDMP Associate) = 26

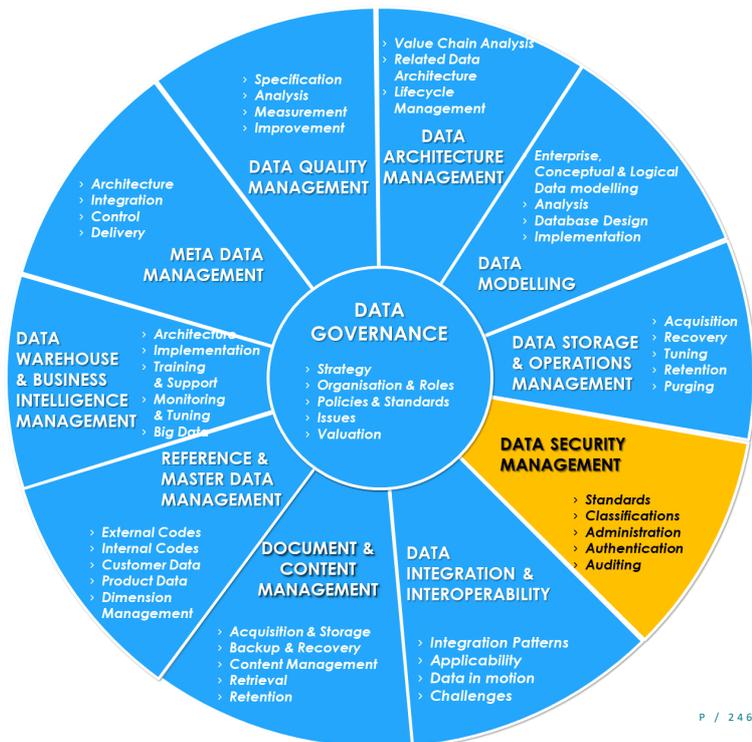
70% (CDMP Practitioner) = 31

80% (CDMP Master) = 35

245

Data Security Management

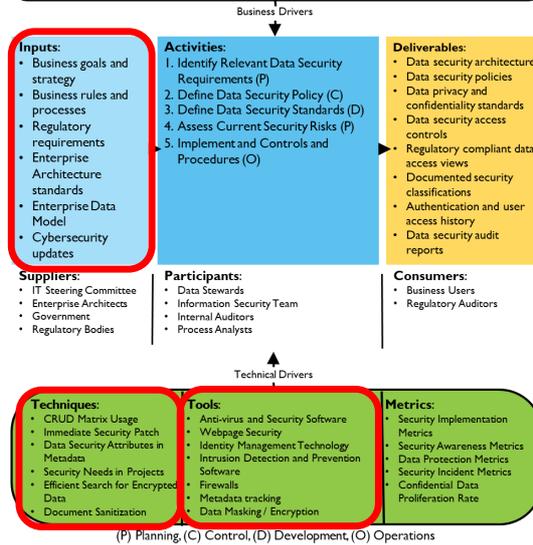
Data Management Process	2%
Big Data	2%
Data Architecture & Lifecycle Management	6%
Document, Records & Content Management	4-5%
Data Ethics	2%
Data Governance	11%
Data Integration & Interoperability	6-7%
Master & Reference Data Management	11%
Data Modelling & Design	11%
Data Quality	11%
Data Security	6%
Data Storage & Operations	4-5%
Data Warehousing & Business Intelligence	11%
Metadata Management	11%



Data Security (DMBoK 2 revised)

Definition: The planning, development, and execution of security policies and procedures to provide proper authentication, authorization, access, and auditing of data and information assets within cultural and regulatory considerations.

Goals:
 1. Enable appropriate, and prevent inappropriate, access to enterprise data assets.
 2. Understand and comply with all relevant regulations and policies for privacy, protection, and confidentiality.
 3. Ensure that the privacy and confidentiality needs of all stakeholders are enforced and audited.



Key Points

- Data Security requirements address “The Four ‘A’s”
 - Authentication Authorization Access Audit
- Data Security is mandated by a range of regulatory provisions across many jurisdictions. It is also basic common sense.
- Breaches of information security can be significantly damaging to the reputation of the organisation, over and above financial or technical damage.
- DMBOK provides a general overview. Prudent Information Managers will refer to relevant ISO standards (ISO27001 / 27701 / 29100) for relevant guidance.



Data Security Guiding Principles



1. Be a responsible trustee of data about all parties
2. Understand and comply with all relevant regulations and guidelines
3. Use CRUD matrices to help map data access needs
4. Ensure Data Security Policy is reviewed and approved by the data governance council (aka DGSC)
5. Identify detailed application security requirements on projects
6. Classify all enterprise data and information products for confidentiality (and other classifications)
7. Set passwords following a set of password complexity guidelines
8. Create security role groups
9. Formally request, track and approve all user and group authorisations
10. Centrally manage user identity data and group membership data
11. Use views or partitions to restrict access to sensitive columns or specific rows
12. Strictly limit and consider every use of shared or service user accounts
13. Monitor data access activity to understand trends



Password Complexity?

<https://www.thehindubusinessline.com/info-tech/cyber-security-log4j-vulnerability-issue-explained/article38061525.ece>



<p>UNCOMMON (NON-GIBBERISH) BASE WORD</p> <p>ORDER UNKNOWN</p> <p>Trøub4dor &3</p> <p>CAPS? COMMON SUBSTITUTIONS NUMERAL PUNCTUATION</p> <p>(YOU CAN ADD A FEW MORE BITS TO FREQUENT FOR THE TREE THAT THIS IS ONLY ONE OF A FEW COMMON FORMATS.)</p>	<p>~28 BITS OF ENTROPY</p> <p>$2^{28} = 3 \text{ DAYS AT } 1000 \text{ GUESSES/SEC}$</p> <p>(PLAUSIBLE ATTACK ON A WEAK REMOTE WEB SERVICE: YES, CRACKING A STOLEN HASH IS FASTER, BUT IT'S NOT WHAT THE AVERAGE USER SHOULD WORRY ABOUT.)</p> <p>DIFFICULTY TO GUESS: EASY</p>	<p>WAS IT TROMBONE? NO, TROUBADOR. AND ONE OF THE O's WAS A ZERO?</p> <p>AND THERE WAS SOME SYMBOL...</p> <p>DIFFICULTY TO REMEMBER: HARD</p>
<p>correct horse battery staple</p> <p>FOUR RANDOM COMMON WORDS</p>	<p>~44 BITS OF ENTROPY</p> <p>$2^{44} = 530 \text{ YEARS AT } 1000 \text{ GUESSES/SEC}$</p> <p>DIFFICULTY TO GUESS: HARD</p>	<p>THAT'S A BATTERY STAPLE.</p> <p>CORRECT!</p> <p>DIFFICULTY TO REMEMBER: YOU'VE ALREADY MEMORIZED IT</p>

THROUGH 20 YEARS OF EFFORT, WE'VE SUCCESSFULLY TRAINED EVERYONE TO USE PASSWORDS THAT ARE HARD FOR HUMANS TO REMEMBER, BUT EASY FOR COMPUTERS TO GUESS.

<https://what3words.com/>

Wordlists. What3words divides the world into a grid of 57 trillion 3-by-3-metre (10 ft × 10 ft) squares, each of which has a three-word address. The company says they do their best to remove homophones and spelling variations; however, at least 32 pairs of English near-homophones still remain.



Founders: Chris Sheldrick; Jack Waley-Cohen; ...

Net income: : -£43.3m (2021);

Operating income: : -£38.4m (2021);

<https://en.wikipedia.org/wiki/What3words> :

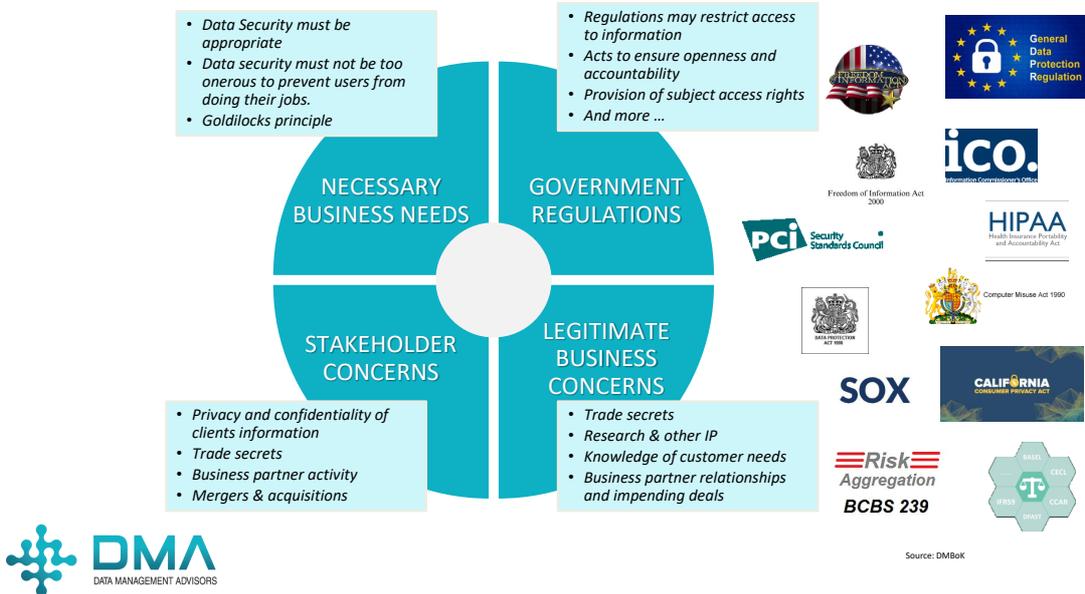
[What3words - Wikipedia](#)



<https://what3words.com/fishery.gala.season>



Sources of Data Security Requirements



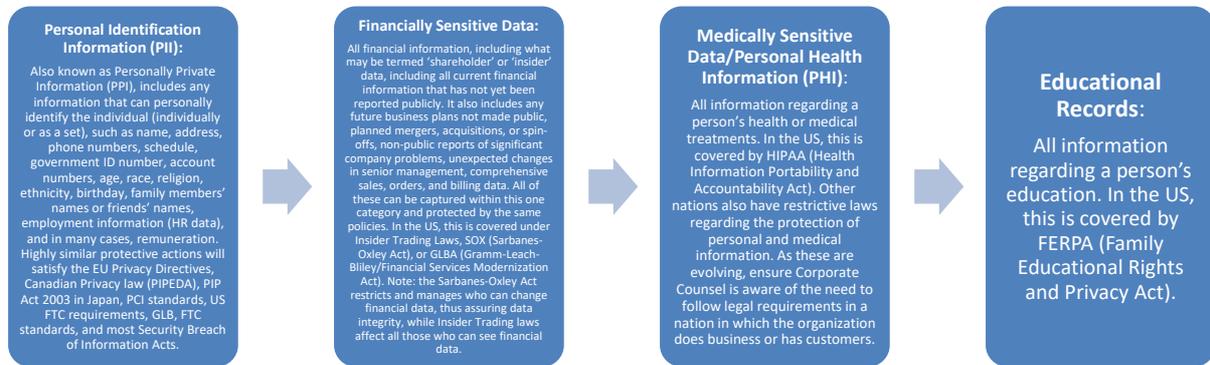
Regulatory and Other Drivers

Security controls are mandated by legislation such as:

- PIPEDA
- PCI-DSS
- EU General Data Protection Regulation (Articles 32-34) GDPR
- HIPAA The Health Insurance Portability and Accountability Act 1996
- Sarbanes-Oxley Act (preventing unauthorised modification of financial transactions)
- BASEL II and Solvency II protecting lineage of data feeding risk models
- CCPA The California Consumer Privacy Act is a state statute intended to enhance privacy rights and consumer protection for residents of California.
- CPRA The California Privacy Rights Act (aka CCPA 2.0)
- BCBS 239 Basel Committee standard number 239 on Banking Supervision. "Principles for effective risk data aggregation and risk reporting" (RDARR).

Increased consumer and media awareness of data security & the “data foot print” means it must be an active consideration for information managers in Business and IT

Regulation Families



.. and more

The California Privacy Rights Act (CPRA) – Effective January 1, 2023

The CPRA applies to for-profit businesses that do business in California and meet any of the following:

- Have a gross annual revenue of over \$25 million;
- Buy, receive, or sell the personal data of 100,000 or more California residents or households; or
- Derive 50% or more of their annual revenue from selling or sharing California residents' personal data.

Virginia Consumer Data Protection Act (CDPA) – Effective January 1, 2023

The CDPA applies to businesses in Virginia, or businesses that produce products or services that are targeted to residents of Virginia, and that:

- During a calendar year, control or process the personal data of at least 100,000 Virginia residents, or
- Control or process personal data of at least 25,000 Virginia residents and derive over 50% of gross revenue from the sale of personal data.

Colorado Privacy Act (CPA) – Effective July 1, 2023

The CPA applies to organizations that conduct business in Colorado or produce or deliver commercial products or services targeted to residents of Colorado and satisfy one of the following thresholds:

- Control or process the personal data of 100,000 Colorado residents or more during a calendar year, or
- Derive revenue or receive a discount on the price of goods or services from the sale of personal data, and process or control the personal data of 25,000 Colorado residents or more.

Connecticut Act Concerning Personal Data Privacy and Online Monitoring (CTPDA) – Effective July 1, 2023

The CTPDA applies to any business that conducts business in the state, or produces a product or service targeted to residents of the state, and meets one of the following thresholds:

- During a calendar year, controls or processes personal data of 100,000 or more Connecticut residents, or
- Derives over 25% of gross revenue from the sale of personal data and controls or processes personal data of 25,000 or more Connecticut residents.

Utah Consumer Privacy Act (UCPA) – Effective December 31, 2023

The UCPA applies to any business that conducts business in the state, or produces a product or service targeted to residents of the state, has annual revenue of \$25,000,000 or more, and meets one of the following thresholds:

- During a calendar year, controls or processes personal data of 100,000 or more Utah residents, or
- Derives over 50% of the gross revenue from the sale of personal data and controls or processes personal data of 25,000 or more Utah residents.



A4



A4



A4



Audit ... Monitor User Authentication and Access Behaviour

- Monitoring authentication and access behaviour is important for:
 - Compliance auditing
 - Alerting security administrators to problems.
- Monitoring helps detect unusual or suspicious transactions.

- **Active**, real-time monitoring is appropriate for systems containing confidential data such as salary information.
- **Passive** monitoring takes snapshots of the system state at regular intervals and compares levels of activity against benchmarks.

Monitoring can be:

- Application specific.
- Implemented for certain users or role-groups.
- Implemented for certain privileges.
- Used to validate data integrity.
- Used to validate configuration.

Defining Standards

- There is no one prescribed way of implementing information security controls in an organisation
- ISO standards and other standards do exist that set out critical requirements for Information Security standards in an organisation:
 - ISO/IEC 29100: Information technology - security techniques - privacy framework Standard. For implementing a privacy framework for use when processing personally identifiable information (PII).
 - ISO/IEC 27001: An international standard on how to manage information security.
 - ISO/IEC 27701: An extension of ISO 27001 to provide the data privacy and information security standards required by GDPR. Implementing ISO 27701 will create a Privacy Information Management System (PIMS)
 - PCI-DSS
- Standards can, and do, influence:
 - Access control
 - Use of mobile devices (e.g. BYOD policy) and associated perimeter security
 - Records and device disposal processes and procedures



<https://www.gov.uk/government/publications/cyber-essentials-scheme-overview>

P / 262

CIA

CONFIDENTIALITY

Preventing the disclosure of information to unauthorised individuals or systems.

INTEGRITY

Preventing the undetectable modification of information.

AVAILABILITY

Ensuring that information is available where and when it is needed.



P / 263

3 categories of controls

Controls are approaches, procedures or other strategies used to address one or more of the information security risks. There are three types (DAMA-I):

ADMINISTRATIVE CONTROLS	<i>“Procedural security”</i> <u>Administrative</u> : Training, Induction, Principles, policies, procedures, standards and guidelines
LOGICAL CONTROLS	<i>“Software security”</i> <u>Logical</u> : Passwords, firewalls, network intrusion detection systems, encryption, access controls lists
PHYSICAL CONTROLS	<i>“Workplace security”</i> <u>Physical</u> : Doors, locks, CCTV, guards, sprinkler system



Risk Assessment:

THREAT <i>An aspect that might be environmental or manmade that has the potential to compromise the confidentiality, integrity or availability of an information asset</i>	VULNERABILITY <i>A weakness that could be exploited to compromise the confidentiality, integrity or availability of an information asset</i>	MITIGATION <i>The activities that can be undertaken, planned or considered to reduce or eliminate the impact of a loss of confidentiality, integrity or availability of an information asset.</i> <i>Auditors want to see all rows in the security risk assessment matrix completed (even if mitigation is unfeasible) in which case “accept” the risk” is included.</i>
PROBABILITY <i>the likelihood that a threat will exploit a vulnerability to compromise the confidentiality, integrity or availability of an information asset</i>	IMPACT <i>A loss of confidentiality, integrity or availability which may result in more significant losses to competitive advantage, revenue, life, property or reputation</i>	

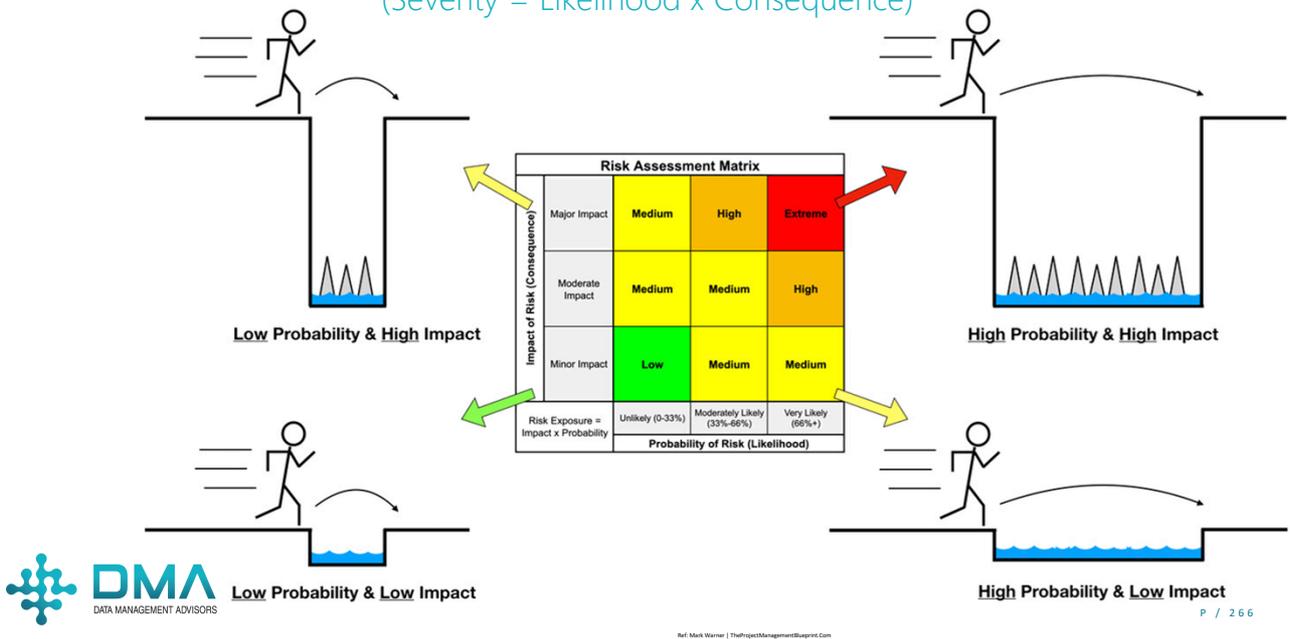
ISO 27001 suggests four ways to treat risks:

1. ‘Terminate’ the risk by eliminating it entirely,
2. ‘treat’ the risk by applying security controls,
3. ‘transfer’ the risk to a third party, or
4. ‘tolerate’ the risk.



Assessment of Risk Exposure = Risk Probability x Impact

(Severity = Likelihood x Consequence)



Risk Classifications



Critical Risk Data (CRD):

Personal information aggressively sought for unauthorized use by both internal and external parties due to its high direct financial value. Compromise of CRD would not only harm individuals, but would result in financial harm to the company from significant penalties, costs to retain customers and employees, as well as harm to brand and reputation.



High Risk Data (HRD):

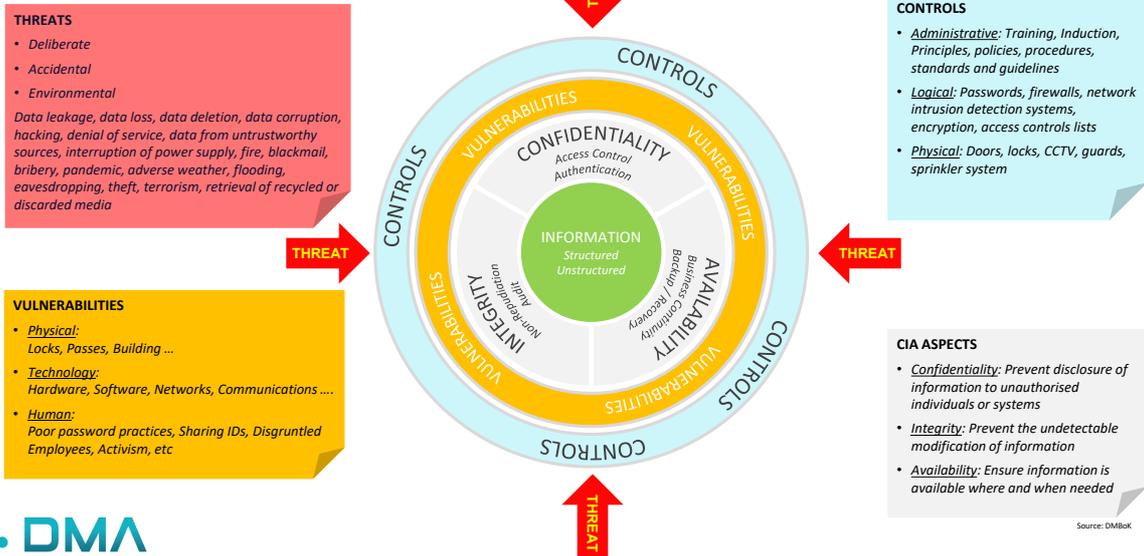
HRD is actively sought for unauthorized use due to its potential direct financial value. HRD provides the company with a competitive edge. If compromised, it could expose the company to financial harm through loss of opportunity. Loss of HRD can cause mistrust leading to the loss of business and may result in legal exposure, regulatory fines and penalties, as well as damage to brand and reputation.



Moderate Risk Data (MRD):

Company information that has little tangible value to unauthorized parties; however, the unauthorized use of this non-public information would likely have a negative effect on the company.

CIA, Threats & Controls in context



IT Security Network Components



Firewall:

A firewall is software and/or hardware that filters network traffic to protect an individual computer or an entire network from unauthorized attempts to access or attack the system. A firewall may scan both incoming and outgoing communications for restricted or regulated information and prevent it from passing without permission (Data Loss Prevention). Some firewalls also restrict access to specific external websites.



Perimeter:

A *perimeter* is the boundary between an organization's environments and exterior systems. Typically, a firewall will be in place between all internal and external environments.



DMZ:

A de-militarized zone (*DMZ*) is an area on the edge or perimeter of an organization, with a firewall between it and the organization. A DMZ environment will always have a firewall between it and the internet. DMZ environments are used to pass or temporarily store data moving between organizations.





- Privilege Escalation**
 Software programs often have bugs that can be exploited. These bugs can be used to gain access to certain resources with higher privileges that can bypass security controls.
- Trojan**
 They masquerade as normal, safe applications, but their mission is to allow a hacker remote access to your computer. In turn, the infected computer can be used as part of a denial-of-service attack and data theft can occur (e.g., keystroke logger).
- Rootkits**
 Rootkits are some of the most difficult to detect. They are activated when your system boots up — before anti-virus software is started. Rootkits allow the installation of files and accounts, or the purposes of intercepting sensitive information.
- Virus**
 A virus is a computer program that, like a medical virus, can replicate and infect other computers.
- Ransomware**
 Ransomware is a subset of malware in which the data on a victim's computer is locked, typically by encryption, and payment is demanded before the ransomed data is decrypted and access returned to the victim..
- Botnets**
 Botnets are created with a Trojan and reside on IRC networks. The bot can launch an IRC client and join chat room in order to spam and launch denial of service attacks.
- Spyware**
 Like Trojans, spyware can pilfer sensitive information, but are often used as advertising tools as well. The intent is to gather a user's information by monitoring Internet activity and transmitting that to an attacker.
- Zero-day**
 A software vulnerability unknown to those who should be interested in its mitigation (including the vendor of the target software). Until the vulnerability is mitigated, hackers can exploit it to adversely affect programs, data, additional computers or a network. An exploit directed at a zero-day is called a zero-day exploit, or zero-day attack.
- Adware**
 Like spyware, adware observes a user's Internet browsing habits. But the purpose is to be able to better target the display of web advertisements.
- Logic bomb**
 They are bits of code added to software that will set off a specific function. Logic bombs are similar to viruses in that they can perform malicious actions like deleting files and corrupting data.
- Spam**
 is unsolicited junk mail. It comes in the form of an advertisement, and in addition to being a time waster, has the ability to consume precious network bandwidth.
- Worm**
 A worm is a specific type of virus. Unlike a typical virus, its goal isn't to alter system files, but to replicate so many times that it consumes hard disk space or memory.



Ref	Question	A	B	C	D	E
DS1	Which of these statements best defines Data Security Management?	The planning, development, and execution of security policies and procedures to provide proper authentication, authorization, access, and auditing of data and information assets	The implementation and execution of checkpoints, checklists, controls, and technical mechanisms to govern the access to information in an enterprise	The definition of controls, technical standards, frameworks, and audit trail capabilities to identify who has or has had access to information	The planning, implementation, and testing of security technologies, authentication mechanisms, and other controls to prevent access to information	Ensuring that Data is Authenticated, Authorised, Accessed and Audited
DS2	Which of these are characteristics of an effective data security policy?	The procedures defined are benchmarked, supported by technology, framework based, and peer reviewed	The defined procedures are tightly defined, with rigid and effective enforcement sanctions, and alignment with technology capabilities	The policies are specific, measurable, achievable, realistic, and technology aligned	The defined procedures ensure that the right people can use and update data in the right way, and that all inappropriate access and update is restricted	The GDPR and FOI regulations are followed
DS3	Apart from security requirements internal to the organisation, what other strategic goals should a Data Security Management system address?	Compliance with ISO29100 and PCI-DSS	Compliance with ISO27001 and HIPAA	Regulatory requirements for privacy and confidentiality AND Privacy and Confidentiality needs of all stakeholders	Ensuring the organisation doesn't engage in SPAM marketing	Ensuring data breaches are reported to the information commissioner in a timely manner
DS4	The implementation and administration of database security is often the responsibility of	The CIO	The Database Administrator	The Database system owner	The Data Governance Council	The System Owner (AKA Custodian)
DS5	What is the role of the Data Governance Council in defining an Information Security policy?	The Data Governance Council should review and approve the high-level Data Security Policy	The Data Governance Council should define the Data Security Policy	The Data Governance Council should implement the Data Security Policy	The Data Governance Council should have no role in Data Security	The Data Governance Council should audit Data Security breaches
DS6	What is the benefit of using role groups to implement data security policies?	It simplifies revoking individual permissions from an individual user	It allows users to be typecast by the administrator	It reduces the amount of effort to assign access rights to users if they inherit rights from their group	It allows for iterative reporting of user access	None, security policies should be set according to individual user needs

AFTER QUIZ 7

1. General (6)
2. Data Quality (9)
3. Data Storage & Operations (3)
4. Master & Reference Data (9)
5. DW / BI + Big Data (7)
6. Data Modelling (9)
7. Data Security (6)

Maximum possible score = 49

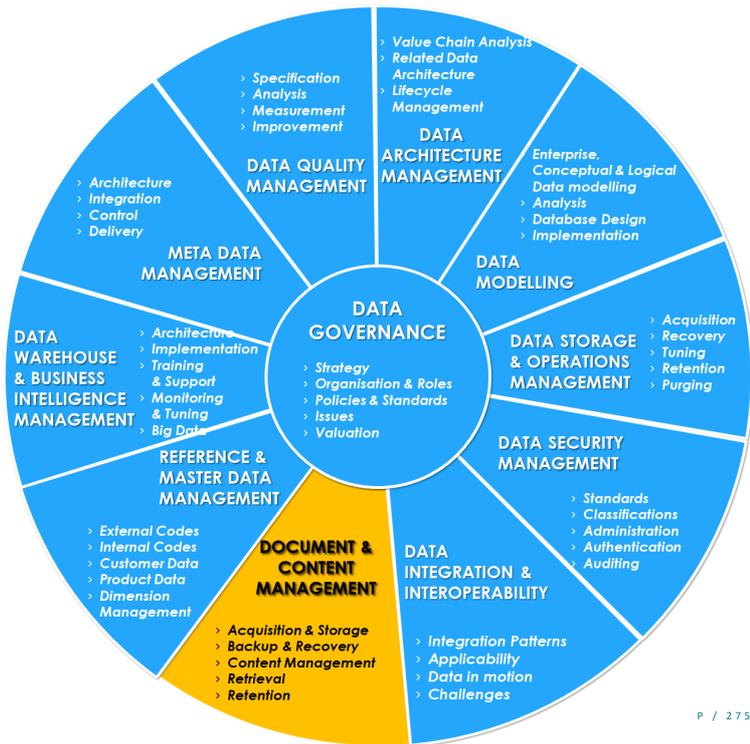
60% (CDMP Associate) = 30

70% (CDMP Practitioner) = 35

80% (CDMP Master) = 40

Document, Records & Content Management

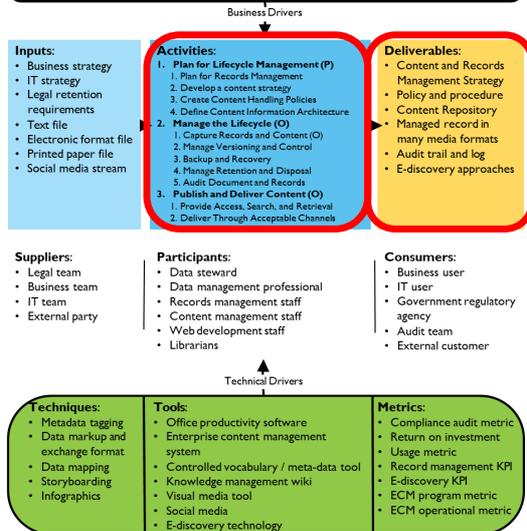
Data Management Process	2%
Big Data	2%
Data Architecture & Lifecycle Management	6%
Document, Records & Content Management	4-5%
Data Ethics	2%
Data Governance	11%
Data Integration & Interoperability	6-7%
Master & Reference Data Management	11%
Data Modelling & Design	11%
Data Quality	11%
Data Security	6%
Data Storage & Operations	4-5%
Data Warehousing & Business Intelligence	11%
Metadata Management	11%



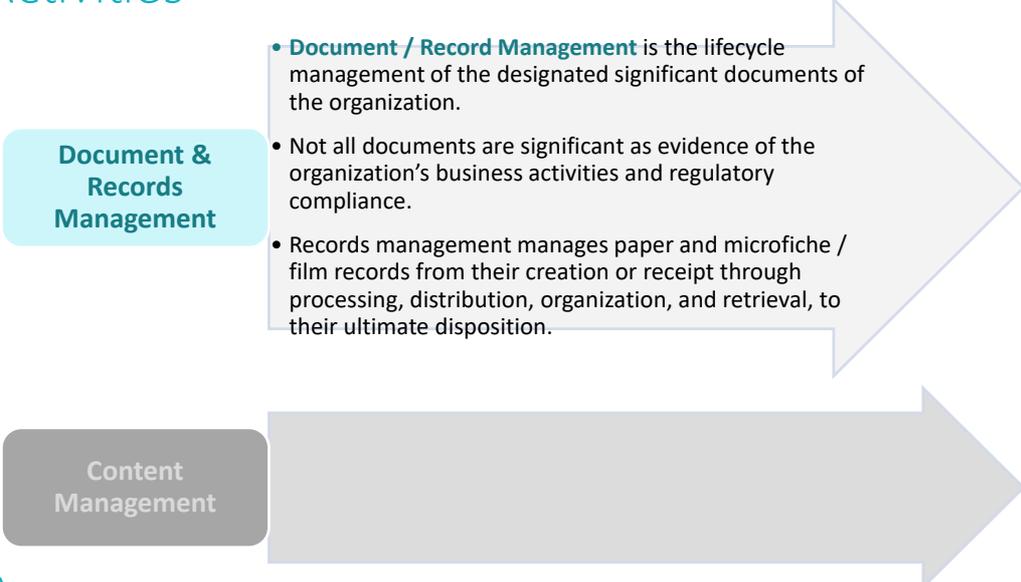
Document and Content Management (DMBoK 2 revised)

Definition: Controlling the capture, storage, access, and use of data and information stored predominantly outside relational databases.

Goals:
 1. To comply with legal obligations and customer expectations regarding Records management.
 2. To ensure effective and efficient storage, retrieval, and use of Documents and Content.
 3. To ensure integration capabilities between structured and unstructured Content.



Main Activities



- **Document / Record Management** is the lifecycle management of the designated significant documents of the organization.
- Not all documents are significant as evidence of the organization's business activities and regulatory compliance.
- Records management manages paper and microfiche / film records from their creation or receipt through processing, distribution, organization, and retrieval, to their ultimate disposition.

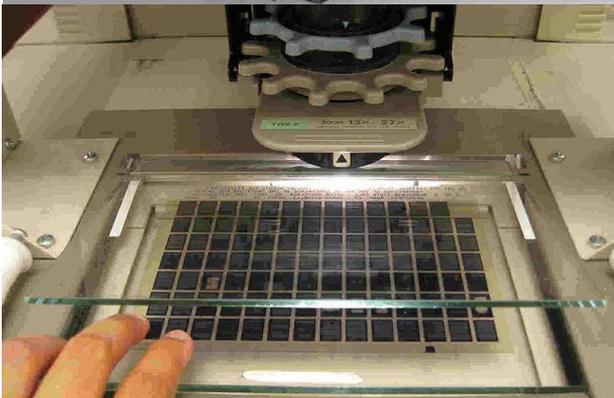
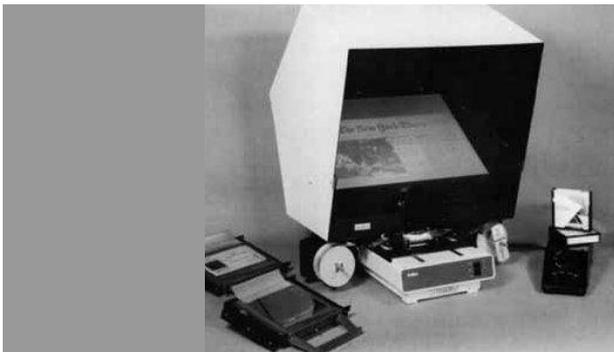


Key Terms

- **Content Management**
The organisation, categorisation, and structure of data / resources so that they can be stored, published and reused in multiple ways.
- **Document Management**
The storage, inventory and control of electronic and paper documents.
- **Records Management**
A subset of Document management. Includes the full lifecycle: from record creation or receipt through processing, distribution, organization, and retrieval, to disposition.

Records can be physical (e.g., documents, memos, contracts, reports or microfiche); electronic (e.g., email content, attachments, and instant messaging); content on a website; documents on all types of media and hardware; and data captured in databases of all kinds.
- **Significant Records / Documents**
Significant as evidence of the organization's business activities and regulatory compliance.
- **Vital Record**
A type of record required to resume an organization's operations the event of a disaster.
- **Taxonomy**
The science or technique of classification.
- **Ontology**
A type of model that represents a set of concepts and their relationships within a domain.





Appliance Category	Class 1		Class 2	
	Visual inspection	PAT test	Visual inspection	PAT test
Fixed	2 yearly	Not required	2 yearly	Not required
Stationary	2 yearly	4 yearly	2 yearly	Not required
IT	2 yearly	4 yearly	2 yearly	Not required
Moveable	Annually	4 yearly	2 yearly	Not required
Portable	Annually	4 yearly	2 yearly	Not required
Cables & Chargers	Annually	4 yearly	2 yearly	Not required
Handheld	6 monthly	Annually	6 monthly	Not required



TAXONOMIES

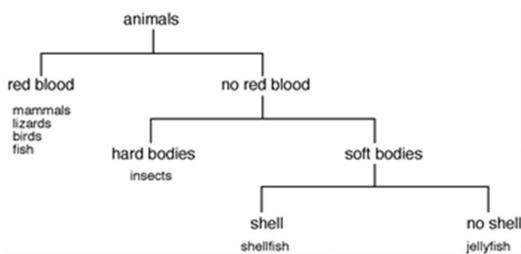
Grouped into five types:

- 1. **Flat Taxonomy** – no relationship among the controlled set of categories (example: list of countries).
- 2. **Hierarchical Taxonomy** – tree structure; for example geography, from continent down to address.
- 3. **A Polyhierarchy** is a tree-like structure with more than one node relation rule. Child nodes may have multiple parents. Those parents may also share grandparents.

- 4. **Facet Taxonomy** – looks like a star where each node is associated with the centre node. Facets are attributes of the object in the centre; for example meta-data where each attribute (creator, title, keywords etc.) is a facet of a content object.
- 5. **Network Taxonomy** – uses both uses both hierarchical and facet structures, for example a thesaurus, or recommender engine (if you liked that, you may also like this...).



Aristotle's classification system (384–322 BC)



Anhaima "bloodless animals" Invertebrata

- Malakia**
"soft bodied"
Fleshy part is exterior and hard part – if any – is interior
Currently: Cephalopoda
- Malakostraca**
"soft-shelled"
Hard part on the outside and fleshy part inside. Hard part has to be crushed rather than scattered
Currently: Malacostraca (in part)
- Ostrakoderma**
"shell skinned"
Hard part outside and fleshy part inside. Hard part can be scattered but not crushed
Currently: Gastropoda, Bivalvia, Echinoidea, Asteroidea, Ascidiacea
- Others, unclassified:**
Currently: Demospongiae, Anthozoa, Polychaeta, Echiura, Amphipoda, Isopoda, Copepoda, Crinoidea, Holothurioidae

Enhaima "blooded animals" Vertebrata

- Ichthyos**
Two peculiarities: fins, because they are essentially swimmers and gills, through which they expel water
Currently: Pisces
- Ketoda**
They have a blow-hole, but no gills
Currently: Cetacea
- Ototoka tetrapoda**
"egg laying tetrapods"
The sea turtle
Currently: Tetrapoda - Reptilia
- Zootoka tetrapoda**
"live-bearing tetrapods"
The seal
Currently: Tetrapoda – Mammalia - Pinnipedia



Example Ontologies

FIBO The Financial Industry Business Ontology (FIBO) is a business conceptual model showing how all financial instruments, business entities and processes work in the financial industry.

CHEMINF The chemical information ontology (cheminf) – contains information entities about chemical entities & includes terms for the descriptors commonly used in cheminformatics software applications and the algorithms which generate them..

FMA - Foundational Model of Anatomy A domain ontology that represents a coherent body of explicit declarative knowledge about human anatomy. Its ontological framework can be applied and extended to all other species.

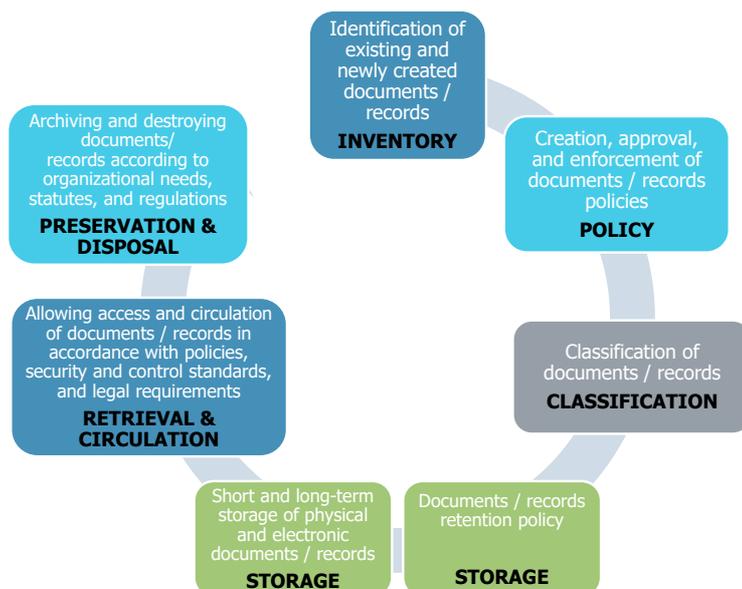
GO - Gene Ontology A collaborative effort to address the need for consistent descriptions of gene products across databases. The GO represents information about biological processes, cellular components and molecular functions.

HPO - Human Phenotype Ontology An ontology that represents phenotypic abnormalities encountered in human disease.

IDO - Infectious Disease Ontology A set of interoperable ontologies covering the infectious disease domain. At the core of the set is a general Infectious Disease Ontology (IDO-Core) of entities relevant to both biomedical and clinical aspects of most infectious diseases. Sub-domain specific extensions of IDO-Core provide coverage of entities relevant to specific pathogens or diseases.



Document / Record Management Lifecycle





Document Control Schemes

- Based on the criticality of the data and the perceived harm that would occur if data were corrupted or otherwise unavailable.
- ANSI Standard 859 ⁽²⁰⁰⁸⁾ describes three levels of control:
 - **Custody control** is the least formal, merely requiring safe storage and a means of retrieval.
 - **Revision control** is more formal, notifying stakeholders and incrementing versions when a change is required. **Requires Administrative Metadata**
 - **Formal control** requires formal change initiation, thorough change evaluation for impact, decision by a change authority, and full status accounting of implementation and validation to stakeholders. **Requires Administrative Metadata**

Document and Records Management Overview

- Lifecycle management of significant organizational documents
- Includes planning, capturing, versioning, retention, and disposal
- Ensures compliance with legal and regulatory requirements
- Supports access, search, retrieval, and audit processes
- Involves collaboration among legal, IT, and business teams



P / 289

Key Principles and Practices in Document Management

- Document control schemes include custody, revision, and formal control
- Classification and taxonomy define document organization and retrieval
- Retention and disposal policies maintain compliance and reduce risk
- Vital records are critical for business continuity post-disaster
- Content management enables reuse and structured publishing



P / 290

Effective Document and Records Management

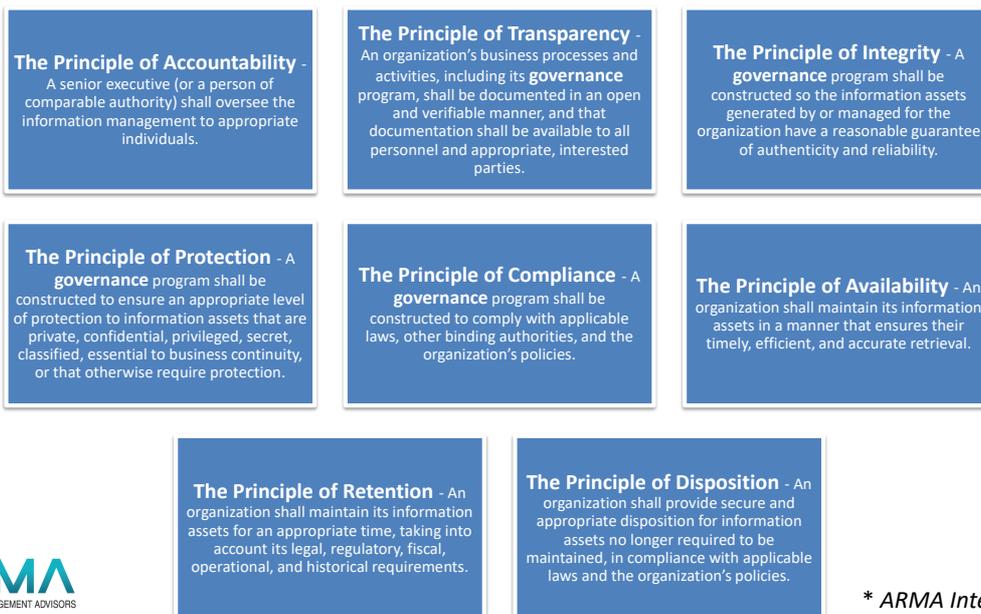
- Clear policies and procedures are essential for retention and disposal.
- Retention policies define maintenance timeframes for valuable documents.
- Inactive documents may be transferred to secondary or off-site storage.
- Oversight ensures privacy, data protection, and identity theft prevention.
- Software versions impact access and readability of electronic records.
- Non-value-added information should be removed to save space and cost.
- Risks exist when retaining records beyond legally required timeframes.
- Challenges include inadequate policies and lack of buy-in for records management.

Many organizations do not prioritize the removal of non-value-added information because:

- Policies are not adequate
- One person’s non-valued-added information is another’s valued information
- Inability to foresee future possible needs for current non-value-added physical and / or electronic records
- There is no buy-in for Records Management
- Inability to decide which records to delete
- Perceived cost of deciding and removing physical and electronic records
- Electronic space is cheap. Buying more space when required is easier than archiving and removal processes



Generally Accepted Recordkeeping Principles GARP*



* ARMA International 292

Main Activities

Document & Records Management



European Commission Article 29 Data Protection Working Party criteria to evaluate anonymization methods.

Anonymity is protected when it is only possible to analyse sizeable "clusters" of individuals who cannot be distinguished from one another based on their attributes

Content Management

- **Content management** is the organization, categorization, and structure of data / resources to be stored, published, and reused in multiple ways.
- **Content** includes data / information, that exists in many forms and in multiple stages of completion within its lifecycle. Content may be found on electronic, paper or other media.
- **The lifecycle of content** can be active, with daily changes through controlled processes for creation, modification, and collaboration of content before dissemination.



Content Delivery Methods

 **Push:** In a push delivery system, users choose the type of content delivered to them on a predetermined schedule. Syndication involves one party creating the content published in many places. Really Simple Syndication (RSS) is an example of a push content delivery mechanism. It distributes content (i.e., a feed) to syndicate news and other web content upon request.

 **Pull:** In a pull delivery system, users pull the content through the Internet. An example of a pull system is when shoppers visit online retail stores.

 **Interactive:** Interactive content delivery methods, such as third-party electronic point of sale (EPOS) apps or customer facing websites (e.g., for enrolment), need to exchange high volumes of real-time data between enterprise applications. Options for sharing data between applications include Enterprise Application Integration (EAI), Changed Data Capture, Data Integration and EII.



Ref	Question	A	B	C	D	E
DOC1	Which of the following are primary deliverables of proper document and record management?	Data from tracking devices, building sensor data	Relational databases, database logs, paper documents	Local drives of laptops, transcripts of phone calls	Spreadsheets, company library books, sales transactions	Managed records in many media formats, e-discovery records, policies and procedures, contracts and financial documents
DOC2	Non-value-added information is often not removed because	The policies are unclear of what is defined as non-value-added so there is no cost driver, and it takes more effort to dispose than to keep.	We might need the information at a later stage	Data is an asset. It is likely to be recognized as valuable in the future.	Legislation is unclear on what should be kept	It should not be removed. All data is value-added
DOC3	When defining your business continuity plan, which of the following should one consider doing?	Have the contracts in place to acquire new hardware in case of technical problems, define policies	Write a report and discuss with management the required budget	Make sure that the data is retained sufficiently long, check that critical data is encrypted, check access rights	Determine the risk, probability and impact, check document backup frequency	Consider written policies and procedures, impact mitigating measures, required recovery time and acceptable amount of disruption, the criticality of the documents



AFTER QUIZ 8

1. General (6)
2. Data Quality (9)
3. Data Storage & Operations (3)
4. Master & Reference Data (9)
5. DW / BI + Big Data (7)
6. Data Modelling (9)
7. Data Security (6)
8. Document, Records & Content Mgt (3)

Maximum possible score = 52

60% (CDMP Associate) = 32

70% (CDMP Practitioner) = 37

80% (CDMP Master) = 42

Data Architecture & Lifecycle Management

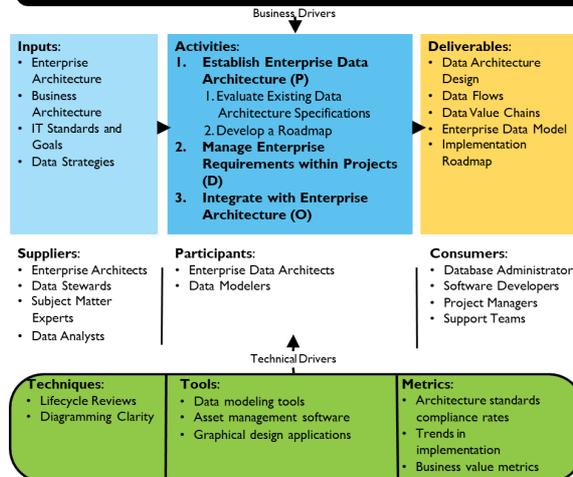
Data Management Process	2%
Big Data	2%
Data Architecture & Lifecycle Management	6%
Document, Records & Content Management	4-5%
Data Ethics	2%
Data Governance	11%
Data Integration & Interoperability	6-7%
Master & Reference Data Management	11%
Data Modelling & Design	11%
Data Quality	11%
Data Security	6%
Data Storage & Operations	4-5%
Data Warehousing & Business Intelligence	11%
Metadata Management	11%



Data Architecture (DMBoK 2 revised)

Definition: Identifying the data needs of the enterprise (regardless of structure) and designing and maintaining the master blueprints to meet those needs. Using master blueprints to guide data integration, control data assets, and align data investments with business strategy.

- Goals:**
1. Identify data storage and processing requirements.
 2. Strategically prepare organizations to quickly evolve their products, services, and data to take advantage of business opportunities inherent in emerging technologies
 3. Design structures and plans to meet the current and long-term data requirements of the enterprise.

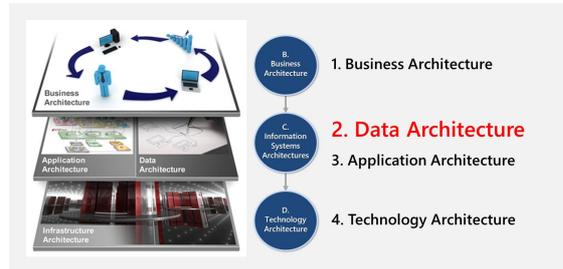


(P) Planning, (C) Control, (D) Development, (O) Operations



Essential Concepts

- Data Architecture appears in all EA Frameworks
- EA Working in practice ...
 - Data model tells you the entity / attribute
 - Process model tells you the action
 - System tells you how the application is programmed to implement the required data processes and handle the data
 - Technology tells you how you are delivering the data management experience to people and other systems
 - Confirmed by the business, application, data and the technology architect



Enterprise Architecture Types & Structures

Enterprise Architecture
Enterprise architecture (EA) is the process of translating business vision and strategy into effective enterprise change by creating, communicating and improving the key requirements, principles and models that describe the enterprise's future state and enable its evolution.

Segment Architecture
Segment architecture is a detailed, formal description of areas within an enterprise, used at the program or portfolio level to organize and align change activity.

Solutions Architecture
Solution architecture is a kind of architecture domain, that aims to address specific problems and requirements, usually through the design of specific information systems or applications.

LEVEL	SCOPE	DETAIL	IMPACT	AUDIENCE
Enterprise Architecture	Agency / Organization	Low	Strategic Outcomes	All Stakeholders
Segment Architecture	Line of Business	Medium	Business Outcomes	Business Owners
Solution Architecture	Function / Process	High	Operational Outcomes	Users and Developers

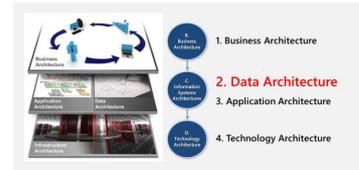
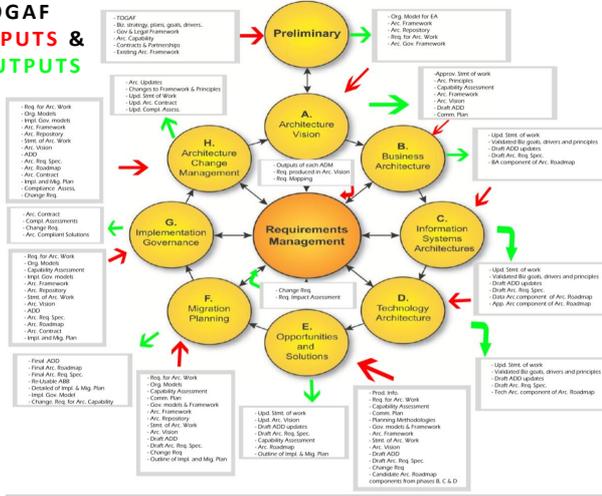


Enterprise Architecture Frameworks

TOGAF

The Open Group Architecture Framework, probably the most widely adopted framework, contains an Architecture Development Method (ADM), content metamodel, and defined artefacts within the business, application, data, and technology domains.

TOGAF INPUTS & OUTPUTS



Enterprise Architecture Frameworks

The Zachman Framework for Enterprise Architecture™

The Enterprise Ontology™

Zachman

The first enterprise architecture framework defines artefacts in a 6 x 6 matrix, with interrogatives (what, how, where, etc.) as columns and stakeholder perspective as rows (executive, business, architect, etc.). It is an ontology, not a methodology for enterprise architecture.

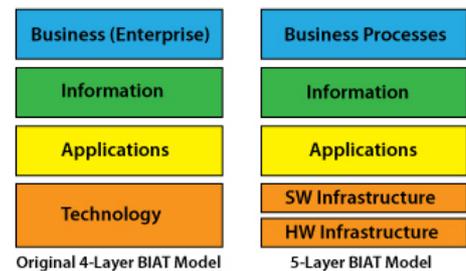
Classification Names	What	How	Where	Who	When	Why	Classification Names
Audience Perspectives	Inventory Identification	Process Identification	Distribution Identification	Responsibility Identification	Timing Identification	Motivation Identification	Model Names
Executive Perspective (Business Context: Finance)	Comp. Integ. → List: Inventory Types	Process Integ. → List: Process Types	Dist. Integ. → List: Distribution Types	Res. Integ. → List: Responsibility Types	Timing Integ. → List: Timing Types	Motiv. Integ. → List: Motivation Types	Scope Contexts (Scope Identification: Entry)
Business Mgmt Perspective (Business Concept: Customer)	Inventory Definition → Business Entity → Business Relationship	Process Definition → Business Transform → Business Input/Output	Distribution Definition → Business Location → Business Connection	Responsibility Definition → Business Role → Business Work Product	Timing Definition → Business Interval → Business Moment	Motivation Definition → Business End → Business Means	Business Concepts (Business Definition: Model)
Architect Perspective (Business Dept: Technology)	Inventory Representation → System Entity → System Relationship	Process Representation → System Transform → System Input/Output	Distribution Representation → System Location → System Connection	Responsibility Representation → System Role → System Work Product	Timing Representation → System Interval → System Moment	Motivation Representation → System End → System Means	System Logic (System Representation: Model)
Engineer Perspective (Business Physics: Business)	Inventory Specification → Technology Entity → Technology Relationship	Process Specification → Technology Transform → Technology Input/Output	Distribution Specification → Technology Location → Technology Connection	Responsibility Specification → Technology Role → Technology Work Product	Timing Specification → Technology Interval → Technology Moment	Motivation Specification → Technology End → Technology Means	Technology Physics (Technology Specification: Model)
Technician Perspective (Business Component: Implementer)	Inventory Configuration → Tool Entity → Tool Relationship	Process Configuration → Tool Transform → Tool Input/Output	Distribution Configuration → Tool Location → Tool Connection	Responsibility Configuration → Tool Role → Tool Work Product	Timing Configuration → Tool Interval → Tool Moment	Motivation Configuration → Tool End → Tool Means	Tool Components (Tool Configuration: Model)
Enterprise Perspective (User: The Enterprise)	Inventory Instantiations → Operations Enables → Operations Relationships	Process Instantiations → Operations Transforms → Operations Inputs/Outputs	Distribution Instantiations → Operations Locations → Operations Connections	Responsibility Instantiations → Operations Roles → Operations Work Products	Timing Instantiations → Operations Intervals → Operations Moments	Motivation Instantiations → Operations Ends → Operations Means	Operations Instances (Implementations: The Enterprise)
Audience Perspectives	Inventory Sets	Process Flows	Distribution Networks	Responsibility Assignments	Timing Cycles	Motivation Intentions	Enterprise Names

© 1997-2011 John A. Zachman. All rights reserved. Zachman® and Zachman International® are registered trademarks of John A. Zachman. To request Permission Use of Copyright, please contact: Zachman.com



Domain	Enterprise Business Architecture	Enterprise Data & Information Architecture	Enterprise Applications Architecture	Enterprise Technology Architecture
Purpose	To identify how an enterprise creates value for customers and other stakeholders	To describe how data should be organized and managed	To describe the structure and functionality of applications in an enterprise	To describe the physical technology needed to enable systems to function and deliver value
Elements	Business models, processes, capabilities, services, events, strategies, vocabulary	Data models, data definitions, data mapping specifications, data flows, structured data APIs	Business systems, software packages, databases	Technical platforms, networks, security, integration tools
Dependencies	Establishes requirements for the other domains	Manages data created and required by business architecture	Acts on specified data according to business requirements	Hosts and executes the application architecture
Roles	Business architects and analysts, business data stewards	Data architects and modelers, business data stewards	Application architects, Solution architects	Infrastructure architects, Solution architects

Enterprise Architecture Domains



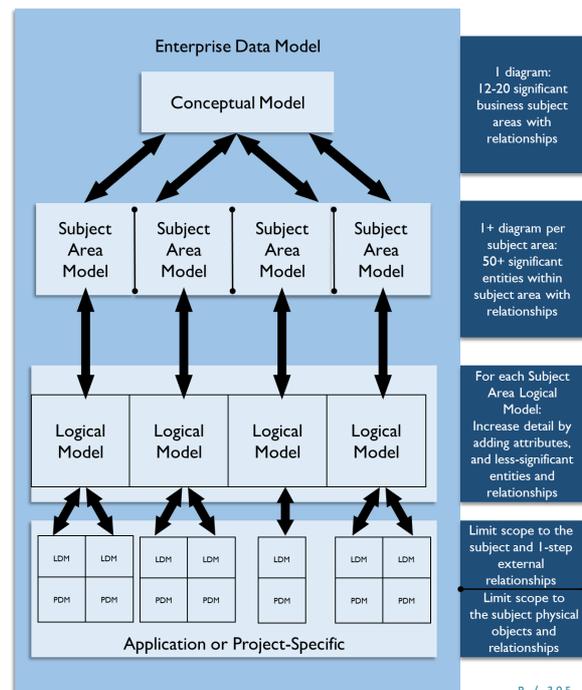
Types of Models in a Typical Data Architecture

Some organizations create an EDM as a stand-alone artifact.

In other organizations, it is composed of data models from different perspectives and at different levels of detail, that consistently describe an organization's understanding of data entities, data attributes, and their relationships across the enterprise.

An EDM includes both universal (Enterprise-wide Conceptual and Logical Models) and application- or project specific data models, along with definitions, specifications, mappings, and business rules.

Adopting an industry standard model (AKA Consensus model) can jumpstart the process of developing an EDM.



Process and Data ARE related

Create
Read
Update
Delete

Major Entities / Data Subject Areas	Business Processes						
	Product development	Marketing & Sales	Industrial preparation	Order management	Manufacturing	Logistics	Invoicing
Product	C	R	U	U	U		
Product Part	C	R	R	U	U		
Manufacturing Plant	U		C	R	R	U	
Customer	R	C		U	R	U	U
Sales Item	C	C	C	U		U	U
Assembly Structure	U		C		U		
Sales Order		U		R	U	U	U
Production Order			U	C	U	U	U
Individual Product					C	R	U
Shipping						C	
Customer's Invoice		U					C



Process and Data ARE related

Create
Read
Update
Delete

Major Entities / Data Subject Areas	Business Processes						
	Product development	Marketing & Sales	Industrial preparation	Order management	Manufacturing	Logistics	Invoicing
Product	C	R	U	U	U		
Product Part	C	R	R	U	U		
Manufacturing Plant	U		C	R	R	U	
Customer	R	C		U	R	U	U
1 Sales Item	C	C	C	U		U	U
Assembly Structure	U		C		U		
2 Sales Order		U		R	U	U	U
Production Order			U	C	U	U	U
Individual Product					C	R	U
3 Shipping						C	
Customer's Invoice		U					C





What is Information Lifecycle Management?

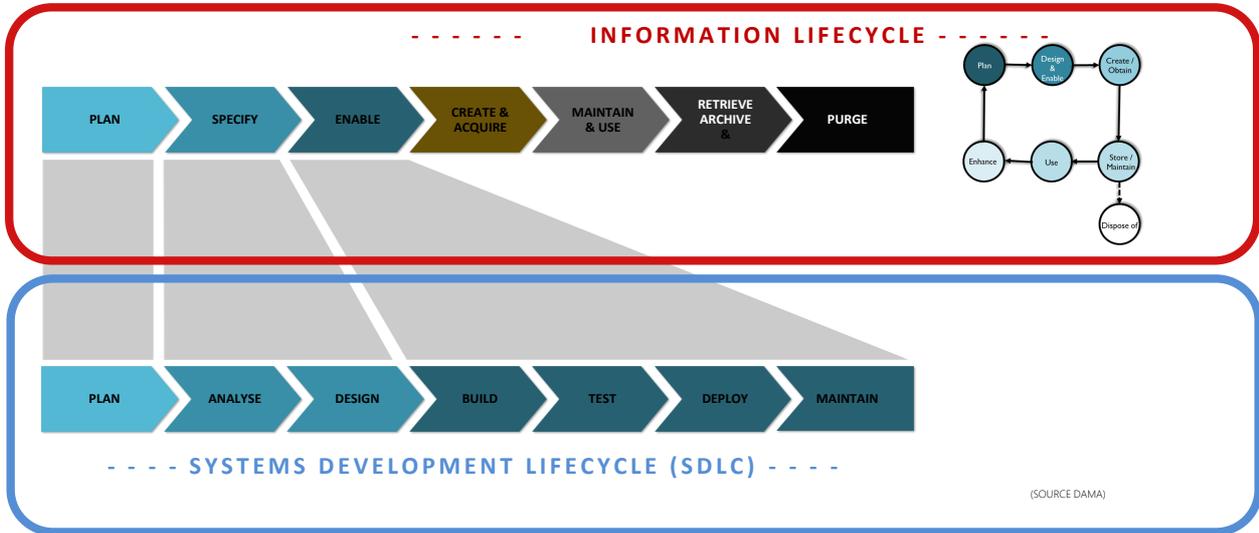
*Firstly, let's understand what we mean by **Information Lifecycle***

The **Information Lifecycle** is the series of phases that information passes through, from inception, through use, to destruction.

Information Lifecycle Management is the means by which best practice is followed at each stage of the **Information Lifecycle**

Information Lifecycle Management defines the planning, implementation and control activities covering the operational, infrastructure, definition and functional areas of data management.

Information Lifecycle & SDLC

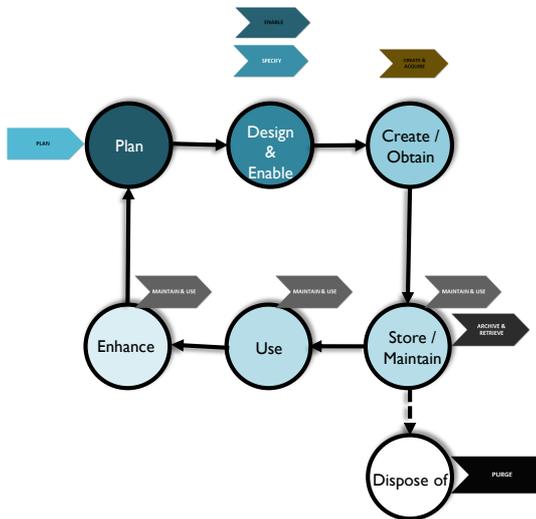


Data is precious

“Data is a precious thing and will last longer than the systems themselves.”



Sir Tim Berners-Lee:
Inventor of the World Wide Web.

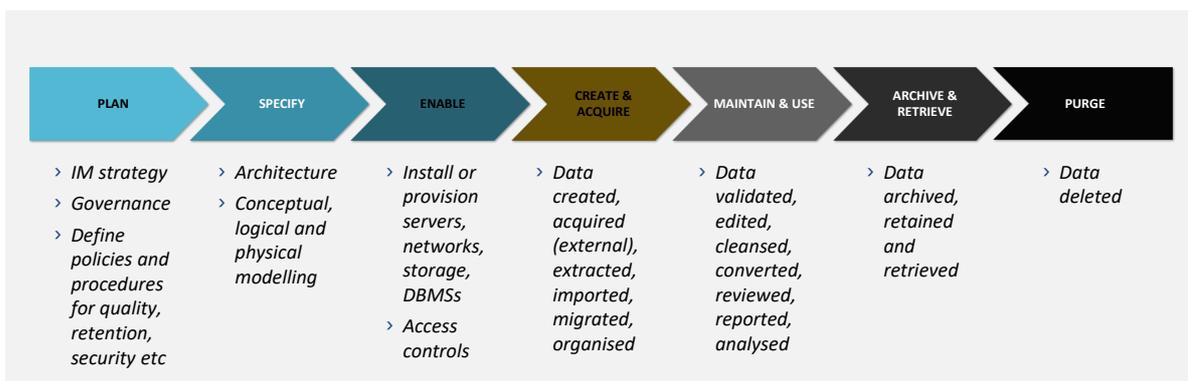


Data Lifecycle 4 "key" activities

1. Data Quality must be managed throughout the data lifecycle
2. Metadata Quality must be managed through the data lifecycle
3. Data Security must be managed throughout the data lifecycle
4. Data Management efforts should focus on the most critical data

What is the Information Lifecycle?

INFORMATION LIFECYCLE (DAMA)



Information Lifecycle Activities

PLAN



- *Define information strategy*
- *Define documentation requirements*
- *Define governance framework*
- *Define monitoring framework*
- *Tool selection*
- *Define information classification framework*
- *Data retention and disposition plan*
- *Staff competencies and training map*
- *Create models for cost of data ownership*

The first stage
of the lifecycle

(SOURCE DAMA)



P / 315

Information Lifecycle Activities

SPECIFY



- *Data requirements gathering*
- *Design data models*
 - *Conceptual*
 - *Logical*
 - *Physical*
 - *As is / To be*
- *Define taxonomies and ontologies*
- *Storage requirements*
- *Define backup requirements*
- *Version control systems*
- *Define warehousing requirements*

Data Modelling
occurs in this
lifecycle stage

(SOURCE DAMA)



P / 316

Information Lifecycle Activities

ENABLE



PROVISION/BUILD

- Networks
- Servers
- Database Management Systems
- Data warehouses
- Embed performance monitoring

APPLY ACCESS CONTROLS

- Based on policies
- Build for availability
- Apply security controls

Infrastructure is provisioned in this lifecycle stage

(SOURCE DAMA)



Information Lifecycle Activities

CREATE & ACQUIRE



The lifecycle stage where information is acquired

- *Treat data as an asset*
- *Buy data from external sources*
 - Postal addresses
 - Dunn and Bradstreet
- *Creating reliable information*
 - Data quality – how good is it?
 - Data currency – is it up-to-date?
 - Status and version control – is it draft or final?
 - Proof of provenance – who created it?
- *Maintain creation information*

(SOURCE DAMA)



Information Lifecycle Activities

MAINTAIN & USE



The lifecycle stage where business value is derived

- *Business value added during use*
 - *Supports business decisions*
 - *Edited and updated*
 - *Data categorised and filed*
 - *Monitor data quality, and fix*
 - *Data may be transformed*
 - *Accessed via data warehousing systems*
 - *Used for reporting*
 - *Data backed-up and restored*
 - *Data distributed and disseminated*
- (SOURCE DAMA)



Information Lifecycle Activities

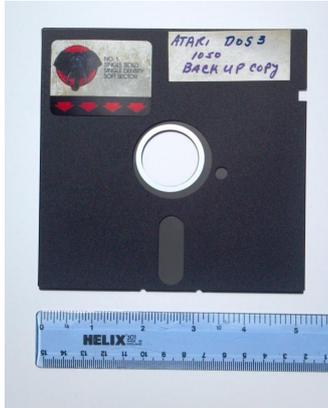
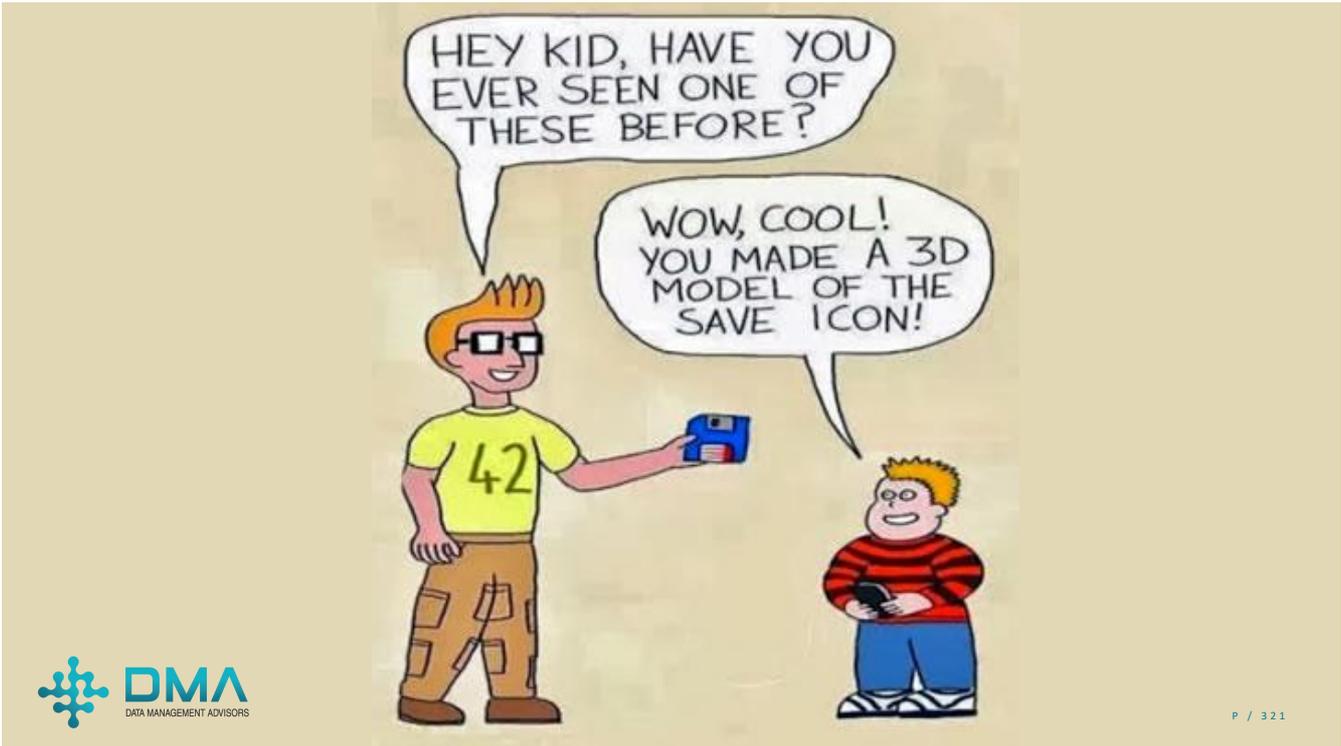
ARCHIVE & RETRIEVE



The lifecycle stage where information is archived (as specified during the planning stage)

- *Policy for archival periods*
 - *How long before data is archived?*
- *Test and verify retrieval methods*
- *Maintain and migrate archive formats*
 - *Media can degrade and become obsolete*
 - *Can old-format physical media still be read? (E.g. floppy disks)*
 - *Can legacy file formats still be read? (E.g. WordPerfect, Lotus 1-2-3)*
 - *Define procedures for validating archives*





Information Lifecycle Activities

PURGE



The lifecycle stage where information is deleted (as specified during the planning stage)

- *According to Retention Plan*
 - Will vary by type of data
- *Consider destruction method*
 - Data is 'deleted' – may still be recoverable
 - Or physical destruction – shred CDs, hard drives, tapes, paper documents
 - "Secure Erase" for magnetic media

(SOURCE DAMA)



P / 3 2 3

Key Points

- The Data Landscape is critical to monitor and govern data asset usage
- The Enterprise Data Model is the fundamental deliverable that:
 - Communicates Business Information Requirements
 - Provides the foundation & guidelines for Data Development
 - Aligns with other business models
 - Provides visual representation of business data
- Data Architecture includes architectural guidelines (reference) for the all the other data management capabilities
 - Operations, DW / BI, Meta-data, Integration, Meta-data, Quality
- The Data Lifecycle is flexible & tailorable for all types of technology



P / 3 2 4

Ref	Question	A	B	C	D	E
DA1	According to the DAMA DMBOK, what parts of the Data Life Cycle are integral parts of the SDLC?	Plan, Specify, Enable	Plan, Create & Acquire, Purge	Specify, Maintain & Use, Purge	Enable, Maintain & Use, Archive & Retrieve	Specify, Enable, Create & Acquire
DA4	Enterprise data architecture defines standard terms for things that are necessary to run an organization. These things are called:	Artefacts	Entities	Taxonomies	Meta-data	Relationships
DA5	The DMBOK identifies which of the following as common stages in the life cycle of the information asset	Plan, Obtain, Store/Share, Maintain, Apply, Dispose	Get, Store, Fix, Use, Purge	Acquire, Integrate, Apply, Share, Dump	Plan, Specify, Enable, Create and Acquire, Maintain & Use, Archive & Retrieve, Purge	Build/Buy, Mix/Merge, Apply, Delete
DA7	The key architecture domains include:	Zachmann, TOGAF, COBIT, and Heath architectures.	business, strategy, application and technology architectures.	business, data, infrastructure and technology architectures.	business, data, application and technology architectures.	process, database, software and technology architectures.



AFTER QUIZ 9

1. General (6)
2. Data Quality (9)
3. Data Storage & Operations (3)
4. Master & Reference Data (9)
5. DW / BI + Big Data (7)
6. Data Modelling (9)
7. Data Security (6)
8. Document, Records & Content Mgt (3)
9. Architecture & lifecycle (4)

Maximum possible score = 56

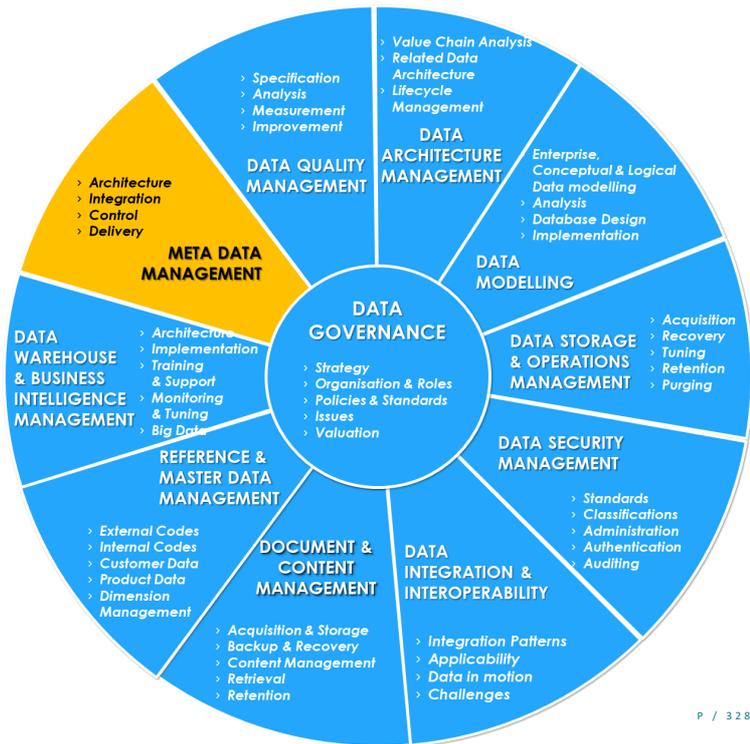
60% (CDMP Associate) = 34

70% (CDMP Practitioner) = 40

80% (CDMP Master) = 45

MetaData Management

Data Management Process	2%
Big Data	2%
Data Architecture & Lifecycle Management	6%
Document, Records & Content Management	4-5%
Data Ethics	2%
Data Governance	11%
Data Integration & Interoperability	6-7%
Master & Reference Data Management	11%
Data Modelling & Design	11%
Data Quality	11%
Data Security	6%
Data Storage & Operations	4-5%
Data Warehousing & Business Intelligence	11%
Metadata Management	11%



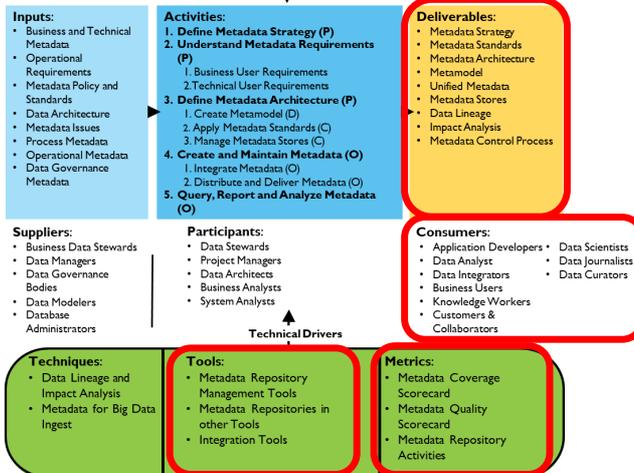
Metadata Management (DMBoK 2 revised)

Definition: Planning, Implementation, and control activities that contributes to the ability to process, maintain, integrate, secure, audit and govern other data

Goals:

1. Provide organizational understanding of business terms and usage including technical lineage.
2. Collect and integrate metadata from diverse sources.
3. Provide a standard way to access metadata and enable known level of trust in data exchange.
4. Ensure metadata quality, consistency, currency, and security.

Business Drivers



(P) Planning, (C) Control, (D) Development, (O) Operations



What is Metadata?



Where do YOU encounter metadata every day?

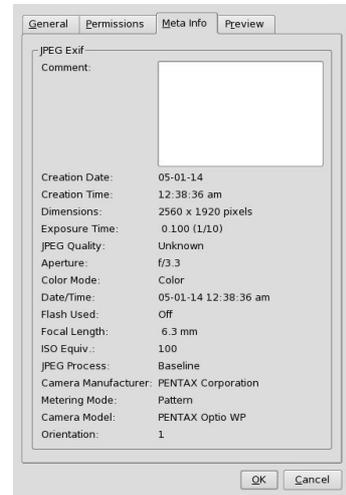


MetaData

DATA



METADATA

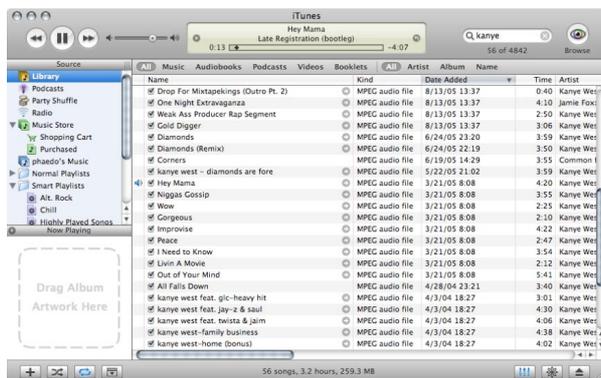


MetaData

DATA



METADATA



Why Do We Need Metadata?



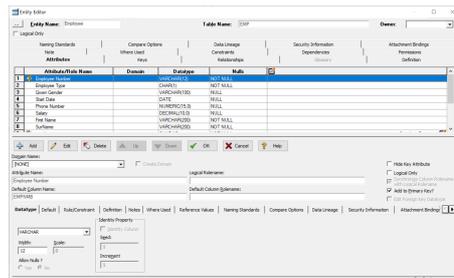
Metadata adds Context and Definition

First Name	Last Name	Company	City	Year Purchased
Kimi	Raikkonen	Komputers R Us	Helsinki	1970
Jenson	Button	The Lord's Store	London	1999
Ayrton	Senna	The Lady's Store	Sao Paulo	1998
Sebastian	Vettel	My Favourite Store	Berlin	2001

Is this the city where the customer lives or the city where the store is located?

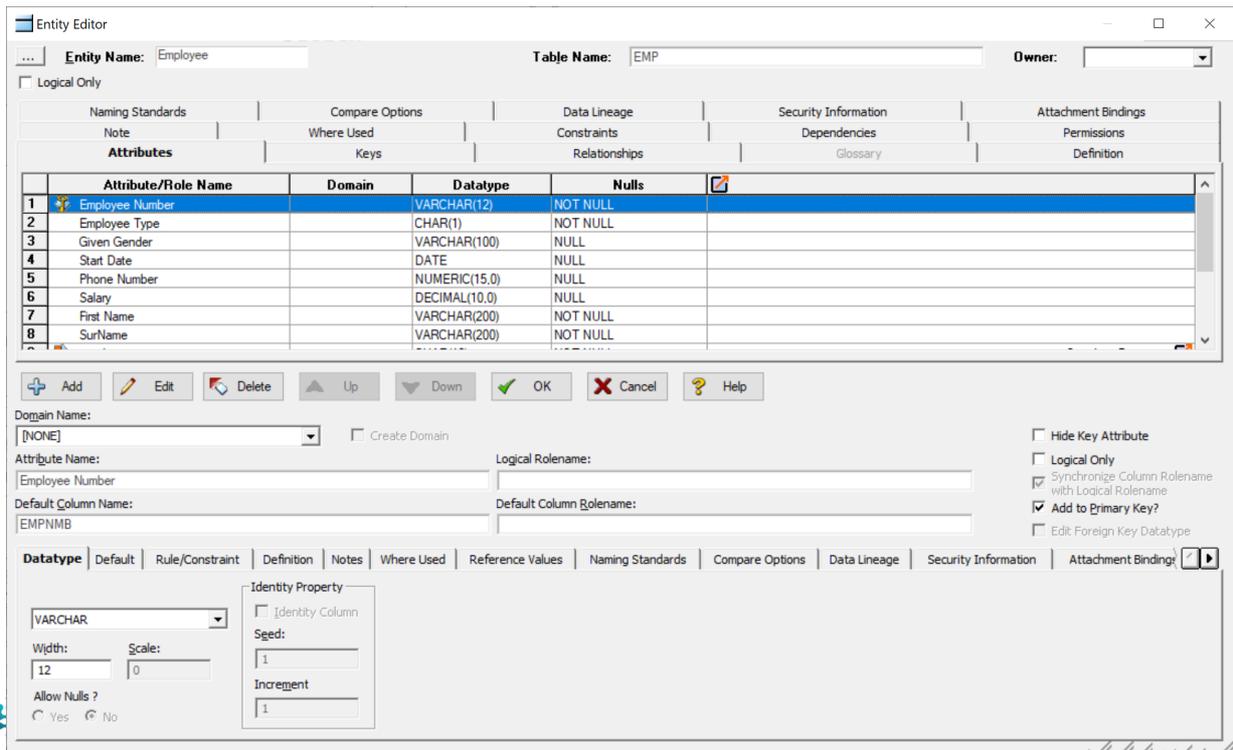
Properties (metadata) of the attributes

Beware – very little (to zero) metadata with big data technologies



Definition	Last Name represents the surname or family name of an individual.
Business Rules	In the Chinese market, family name is listed first in salutations.
Format	VARCHAR(30)
Abbreviation	LNAME
Required	YES
Etc.	Numerous technical & business metadata including security, privacy, nullability, primary key, etc.





Metadata covers the 6 interrogatives of Data

Who	What	Where	Why	When	How
Who is regulating or auditing this data?	What are the technical naming standards for database implementation?	Are there local privacy or security policies that regulate this data?	Why are we regulating or auditing this data?	When do regulators and / or auditors require the data?	How are the auditors / regulators permitted to request this data?
Who created this data?	What is the business definition of this data element?	Where is this data stored?	Why are we storing this data?	When was this data created?	How is this data formatted? (character, numeric, etc.)
Who is the Steward of this data?	What are the business rules for this data?	Where did this data come from?	Why are we storing this data ie its usage & purpose?	When was this data last updated?	How many databases or data sources store this data?
Who is using this data?	What is the security level or privacy level of this data?	Where is this data used & shared?	Why is this data managed, i.e. the business motivations?	When does it need to be purged / deleted?	How / by what mechanisms / devices is this data handled?
Who "owns" this data?	What is the abbreviation or acronym for this data element?	Where is the backup for this data?	Why is the data classified as PII?	When should the data be updated?	How is the data provided to the regulators?
Who is permitted to use this data?	What are the business drivers for using this data?	Where is the disaster recovery location for this data?	Why is this data business critical?	When was the data last quality checked?	How long should it be stored?

Types & typical Sources of MetaData

- Application Code
- Big Data platforms (but not much here)
- Business Intelligence (BI) Tools
- Business Process Models
- CMDB
- COBOL copybooks
- Data Models (extended metadata)
- Data Transformation (e.g. ETL Tools)
- Data Quality Tools
- ERP, CRM, and Package Applications
- Humans
- Internet of Things (IoT)
- Open Data
- Photos / Images (EXIF)
- PREMIS (OAIS Reference Model (ISO 14721) Open Archival Information System (or OAIS) PREservation Metadata: Implementation Strategies (PREMIS))
- Relational databases (system catalogue)
- Social Media
- Spreadsheets
- Text Documents
- Transaction logs
- XML
- and more



P / 339

NoSQL – Key Value Databases

- NoSQL Databases are often optimal solutions for flexibility & performance in certain scenarios.

- One common NoSQL database is a key-value pair database (e.g. Redis, RIAK, Oracle NoSQL, etc.)
- They can support extremely high volumes of records & state changes per second through distributed processing and distributed storage.
- Use cases include: Managing user sessions in web applications, online gaming, online shopping carts, etc.

4 types of NoSQL databases:

1. Document
2. Key Value
3. Graph
4. Columnar

Key	Value
1839047	John Doe, Prepaid, 40.00
9287320	01/01/2008, 50.00, Green

- The structure is often created by the **application code**, not within a database or metadata structure.

- Metadata for NoSQL databases is typically minimal or non-existent.



- The structure & metadata is generally determined by the application code

P / 340

Types & typical Sources of MetaData

The OCCURS clause marks an array.
 The DEPENDING ON clause marks a counter field for the array, if one exists.
 The array ASSIGNMENTS is a nested array within COURSES.

```

01 STUDENT.
20 ID
20 FIRST_NAME
20 LAST_NAME
*
20 DATE_OF_BIRTH
20 NUMOF_COURSES
20 NUMOF_BOOKS
20 COURSES.
  25 COURSE OCCURS 8 TIMES DEPENDING ON NUMOF_COURSES.
    30 COURSE_ID
    30 COURSE_TITLE
    30 INSTRUCTOR_ID
    30 NUMOF_ASSIGNMENTS
    40 ASSIGNMENT_TYPE
    40 ASSIGNMENT_TITLE
    40 DUE_DATE
    40 GRADE
  25 BOOKS.
    25 BOOK OCCURS 1 TO 5 TIMES DEPENDING ON NUMOF_BOOKS.
      30 ISBN
      30 RETURN_DATE
    PTC 9(8).
    PTC X(32).
    PTC X(32).
    PTC S9(8) COMP.
    PTC S9(8) COMP.
    PTC 9(4) COMP.
    PTC 9(4) COMP.
    PTC 9(4) COMP.
    PTC 9(8).
    PTC X(48).
    PTC 9(8).
    PTC X(4) COMP.
    PTC X(12).
    PTC X(48).
    PTC S9(8) COMP.
    PTC S9(8) COMP.
    PTC X(10).
    PTC X(8) COMP.
    
```

```

BookCatalog.xsd
<?xml version="1.0" encoding="UTF-8"?>
<xsd:schema xmlns:xsd="http://www.w3.org/2001/XMLSchema"?>
  <xsd:element name="BookCatalog">
    <xsd:complexType>
      <xsd:sequence>
        <xsd:element ref="Book" minOccurs="1"
          maxOccurs="unbounded"/>
      </xsd:sequence>
    </xsd:complexType>
  </xsd:element>
  <xsd:element name="Book">
    <xsd:complexType>
      <xsd:sequence>
        <xsd:element ref="Title" minOccurs="1" maxOccurs="1"/>
        <xsd:element ref="Author" minOccurs="1" maxOccurs="1"/>
        <xsd:element ref="Date" minOccurs="1" maxOccurs="1"/>
        <xsd:element ref="ISBN" minOccurs="1" maxOccurs="1"/>
        <xsd:element ref="Publisher" minOccurs="1" maxOccurs="1"/>
      </xsd:sequence>
    </xsd:complexType>
  </xsd:element>
  <xsd:element name="Title" type="xsd:string"/>
  <xsd:element name="Author" type="xsd:string"/>
  <xsd:element name="Date" type="xsd:string"/>
  <xsd:element name="ISBN" type="xsd:string"/>
  <xsd:element name="Publisher" type="xsd:string"/>
</xsd:schema>

```



Social Media Metadata

Metadata from social media, such as Twitter, can help identify trend and sentiment analysis.

- Author
- Text/Content of Tweet
- Hash Tag
- # of Retweets



- Date/Time.
- Location.
- Language.
- Device / OS.
- User ID.
- ID of Followers.
- ID of users who "liked" tweet.
- ID of users retweeting tweet.
- Spaces.
- Direct Messages.
- Lists.
- Trends.
- Media.
- Places.
- Etc.



Types of metadata 3 CATEGORIES - DMBok



Categories in relation to where it originates vs where it is used

Business metadata

Relates business perspective to the metadata user (e.g., business data definitions, regulatory or contractual constraints, Data Quality statements).

Examples include:

- Definitions and descriptions of data sets, tables, and columns
- Business rules, transformation rules, calculations, and derivations
- Data models
- Data quality rules and measurement results
- Schedules by which data is updated
- Data provenance and data lineage
- Data standards
- Designations of the system of record for data elements
- Valid value constraints
- Stakeholder contact information (e.g., data owners, data stewards)
- Security/privacy level of data
- Known issues with data
- Data usage notes

Technical

The technical details of the data in IT systems (e.g., systems that store data, processes that move and use it)

Examples include:

- Physical database table and column names
- Column properties
- Database object properties
- Access permissions
- Data CRUD (create, read, update and delete) rules
- Physical data models, including data table names, keys, and indexes
- Documented relationships between the data models and the physical assets
- Data Integration / loading job details
- File format schema definitions (DDL)
- Source-to-target mapping documentation
- Data lineage documentation, including upstream and downstream change impact information
- Program and application names and descriptions
- Content update cycle job schedules and dependencies
- Recovery and backup rules
- Data access rights, groups, roles

Operational metadata

Targeted at IT operations users' needs (e.g., data archiving and retention rules, audit rules, purge criteria)

Examples include:

- Logs of job execution for batch programs
- History of extracts and results
- Schedule anomalies
- Results of audit, balance, control measurements
- Error Logs
- Reports and query access patterns, frequency, and execution time
- Patches and Version maintenance plan and execution, current patching level
- Backup, retention, date created, disaster recovery provisions
- SLA requirements and provisions
- Volume and usage patterns
- Data archiving and retention rules, related archives
- Purge criteria
- Data sharing rules and agreements
- Technical roles and responsibilities, contacts



Types of metadata – Others (Real World)



• Process metadata

Other system elements e.g.

- Process name,
- Data Stores & Data Involved,
- Government/ Regulatory Bodies,
- Roles & Responsibilities,
- Process Handoffs,
- Process Dependencies
- Decomposition ...

• Data Stewardship / Governance metadata

Data about stewards and stewardship processes (e.g.

- Data owners,
- Data subject areas,
- Data users, Data stewards,
- Data Classification ...



Types of metadata – Information Science



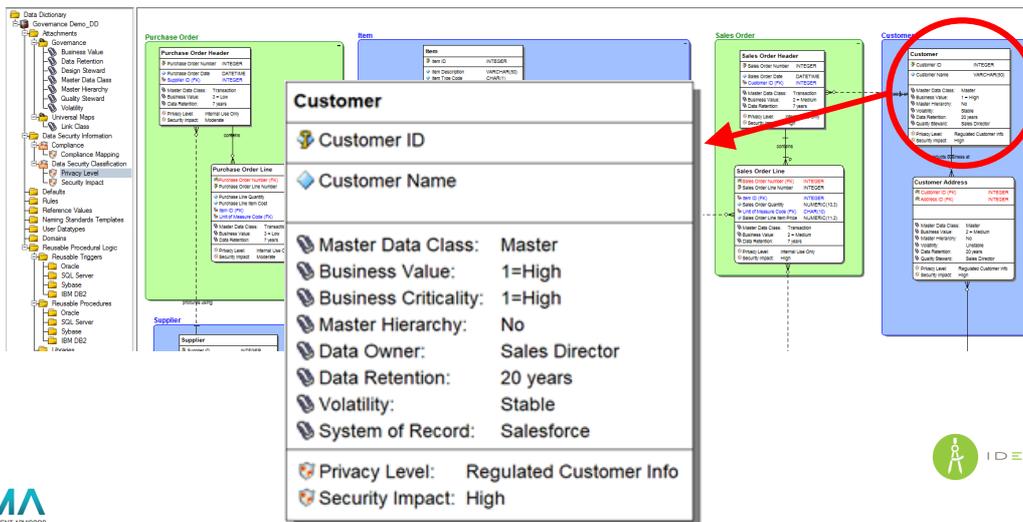
Separate to IT i.e. information science, Metadata is described using 3 different categories:

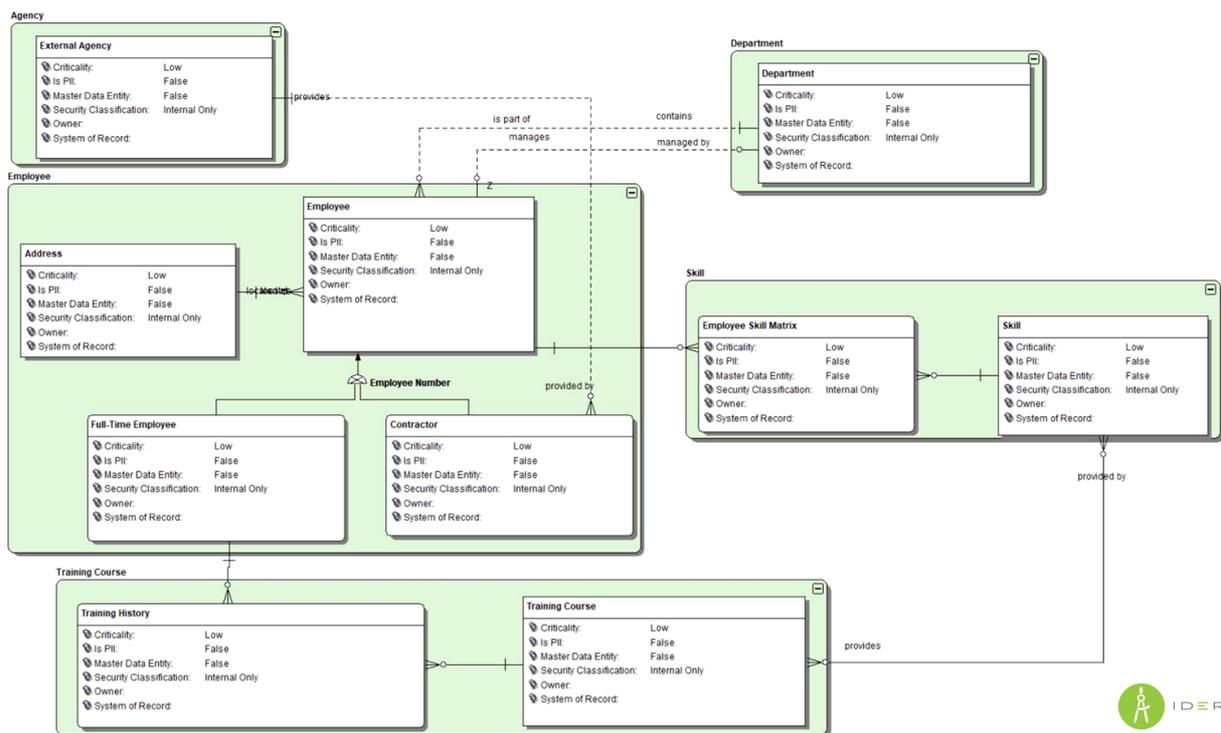
- Descriptive Metadata**
 e.g., title, author, reference #, subject, keywords. This describes a resource and enables identification and retrieval.
- Structural Metadata**
 This describes relationships within and among resources and their component parts e.g., number of pages, number of chapters, index, references, etc.
- Administrative Metadata**
 e.g., version numbers, archive dates, last updated, approvers etc. This is used to manage resources over their lifecycle – essential for some Document Control Schemes (DCS)



Data Governance & Models – Metadata extensions

(ER/Studio Attachments)





Metadata & Data Governance

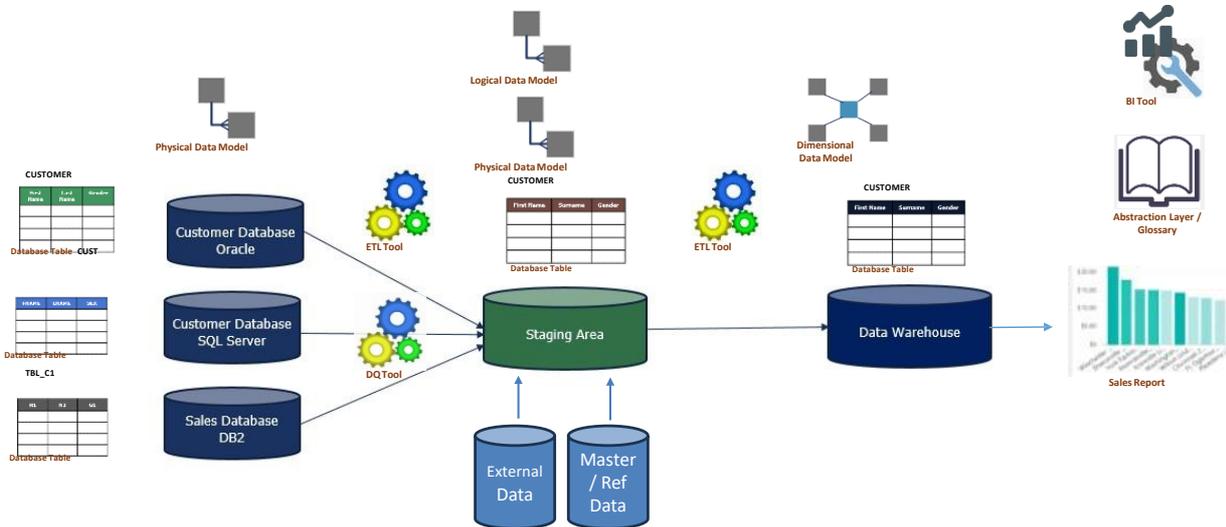
Major data governance aspects held in metadata include:

- **Business Glossary:** Defining the business metadata definitions for business data elements.
- **Data Dictionary:** Defining the technical metadata definitions for data(base)objects (fields, files, tables etc)
- **Data Stewardship:** Aligning data stewardship or ownership roles to key data objects.
- **Data Standards:** Metadata standards provide governance and rules for current and future development, based on both business and IT input.
- **Privacy & Security:** Identification of privacy and security levels for business data is a key aspect of data governance can be managed with metadata definitions.
- **Traceability & Audit:** Understanding of how data is used across the organization, how key financial figures are calculated, etc.

Metadata is a key enabler for a Data Governance initiative

Data Lineage

Data Warehousing Example Metadata for CUSTOMER exists in a number tools & data stores



Data Values Standards (some examples)

COUNTRY Codes:

- **ISO 3166**
 - Alpha-2, Alpha-3, Numeric codes for countries, codes for subdivisions and formerly used codes
- **FIPS 10**
- **IOC**
- **Licence Plate**
- **Internet Domain (ICANN)**

CURRENCY Codes:

- **ISO 4217**
 - ISO 3166 Alpha-2 + 1 char Currency
 - E.g. GBP, USD, CHF

DATE / TIME:

- **ISO 8601**
 - D = yyyyymmdd
 - T = hhmmss

LANGUAGE Codes:

- **ISO 639**
 - Alpha-2, Alpha-3, Name, Scope, Type
 - EN, ENG, English, Individual, Living
 - LA, LAT, Latin, Individual, Ancient

UNIT OF MEASUREMENT (UOM):

- **7 SI base units**
 - Length - meter (m)
 - Time - second (s)
 - Amount of substance - mole (mole)
 - Electric current - ampere (A)
 - Temperature - kelvin (K)
 - Luminous intensity - candela (cd)
 - Mass - kilogram (kg)

CLASSIFICATION OF DISEASES

- World Health Organisation (WHO) International Classification of Diseases (ICD). The standard for all clinical & research purposes.

Country	FIPS10-4	ISO 3166
Aland Islands	–	AX
Algeria	AG	DZ
American Samoa	AQ	AS
Antigua & Barbuda	AC	AG
Australia	AS	AU
Austria	AU	AT
Bahrain	BA	BH
Bosnia and Herzegovina	BK	BA
Gabon	GB	GA
Saudi Arabia	SA	SA
South Sudan	-	SS
Sudan	SU	SD
United Kingdom of Great Britain and Northern Ireland	UK	GB



Data Values Standards (some examples)

DATES: ISO 8601

Dates constructs in ISO 8601	
Date	2021-10-23
Date and time in UTC	2021-10-23T15:08:07+00:00 2021-10-23T15:08:07Z 20211023T150807Z
Week	2021-W42
Week with weekday	2021-W42-6
Date without year	--10-23
Ordinal date	2021-296



Exchangeable image file format (officially **Exif**, according to JEIDA/JEITA/CIPA specifications) is a standard that specifies the [formats](#) for [images](#), [sound](#), and ancillary tags used by [digital cameras](#) (including [smartphones](#)), [scanners](#) and other systems handling image and sound files recorded by digital cameras. The specification uses the following existing file formats with the addition of specific [metadata](#) tags: [JPEG discrete cosine transform](#) (DCT)^[2] for compressed image files, [TIFF](#) Rev. 6.0 (RGB or YCbCr) for uncompressed image files, and [RIFF WAV](#) for audio files (Linear [PCM](#) or ITU-T [G.711](#) μ -Law PCM for uncompressed audio data, and [IMA-ADPCM](#) for compressed audio data).^[3] It is not used in [JPEG 2000](#) or [GIF](#). This standard consists of the Exif image file specification + Exif audio file specification.

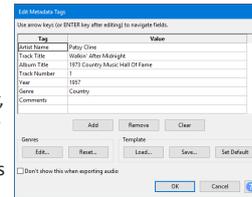


Extended Display Identification Data (EDID)

A [metadata](#) format for [display devices](#) to describe their capabilities to a video source (e.g. [graphics card](#) or [set-top box](#)). The data format is defined by a standard published by the Video Electronics Standards Association ([VESA](#)). The EDID data structure includes manufacturer name and serial number, product type, [phosphor](#) or [filter](#) type, timings supported by the display, display size, [luminance](#) data and (for digital displays only) [pixel](#) mapping data. [DisplayID](#) is a VESA standard targeted to replace EDID and [E-EDID](#) extensions with a uniform format suited for both PC monitor and consumer electronics devices



ID3 is a [metadata](#) container most often used in conjunction with the [MP3 audio file format](#). It allows information such as the title, artist, album, track number, and other information about the file to be stored in the file itself. There are two unrelated versions of ID3: ID3v1 and ID3v2. ID3v1



takes the form of a 128-[byte](#) segment at the end of an MP3 file containing a fixed set of data fields. ID3v1.1 is a slight modification which adds a "track number" field at the expense of a slight shortening of the "comment" field. ID3v2 is structurally very different from ID3v1, consisting of an extensible set of "frames" located at the start of the file, each with a frame identifier (a three- or four-byte string) and one piece of data. 83 types of frames are declared in the ID3v2.4 specification, and applications can also define their own types. There are standard frames for containing cover art, BPM, copyright and license, lyrics, and arbitrary text and URL data.



Industry Metadata Standards

- **OMG (Common Warehouse Metadata (CWM), Information Management Metamodel (IMM), MDC Open Information Model (OIM), XML, UML, SQL)**
- **World Wide Web Consortium (W3C): RDF (Relational Definition Framework)**
- **Dublin Core: Dublin Core Metadata Initiative (DCMI)**
- **Distributed Management Task Force (DTMF): Web-Based Enterprise Management (WBEM)**
- **XML Metadata Interchange (XMI)**
- **Metadata Standards for Unstructured Data**
- **ISO / IEC 11179 Metadata Registry Standard for naming of concept, data element concept, conceptual domain, data element, and value domain.**

CWM defines a specification for modelling metadata for relational, non-relational, multi-dimensional, and most other objects found in a data warehousing environment. Uses UML, MOF & XMI

Dublin Core Schema is a small set (15) of vocabulary terms that can be used to describe digital resources (video, images, web pages, etc.), as well as physical resources such as books or CDs and objects like artworks.

OMG standard for exchanging metadata information via XML. Current version 2.5.1

ISO metadata registries (MDR) standard



- **Data Catalog Vocabulary (DCAT) from W3C**

Meta-Object Facility (MOF) An Object Management Group standard for model-driven engineering. Its purpose is to provide a type system for entities in the CORBA architecture and a set of interfaces through which those types can be created and manipulated.

Metadata & EDI

For EDI (Electronic Data Interchange), the primary metadata standard is UN/EDIFACT, an international standard developed by the United Nations, widespread outside of North America; Within North America, the most common EDI metadata standard is ANSI ASC X12 (X12).

- **Data element structure:**

EDI standards define the structure of data elements within a message, including data types, codes, and lengths, ensuring consistency in data exchange.

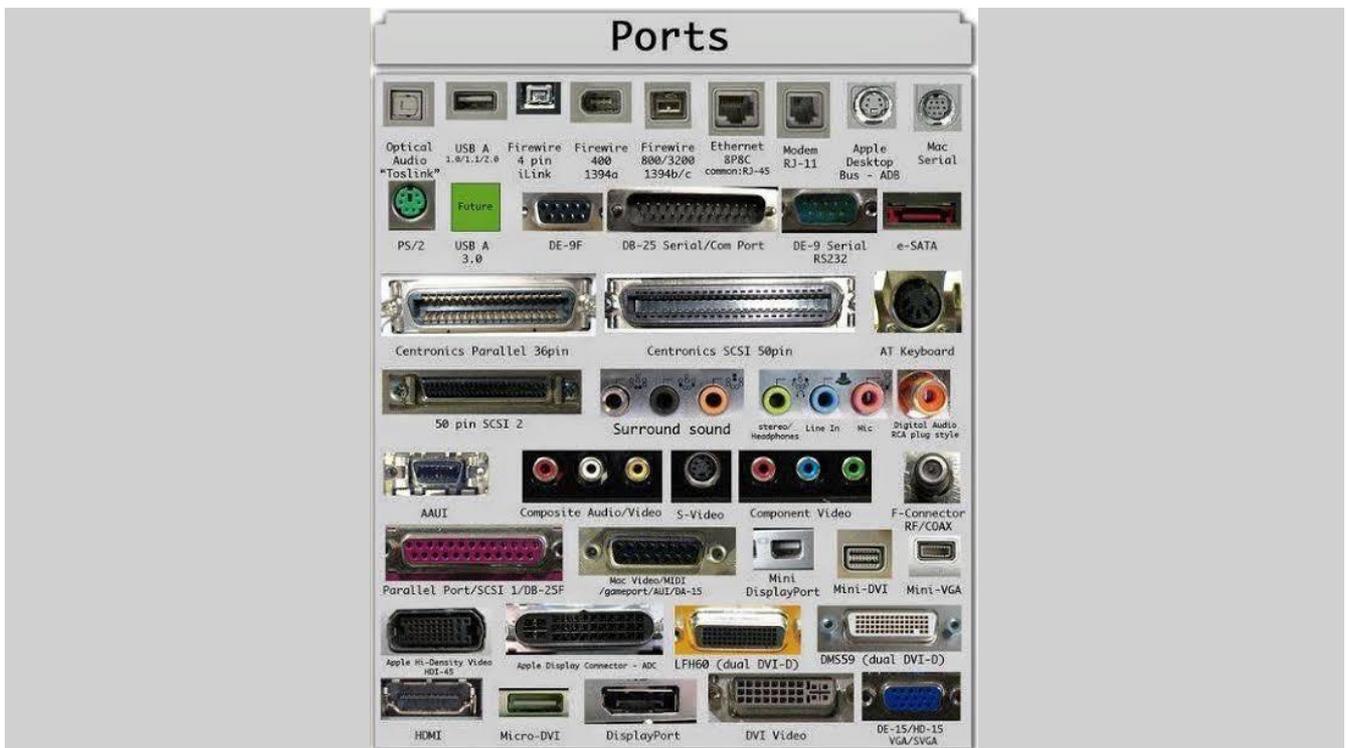
- **Message types:**

Each EDI standard includes various message types, such as purchase orders, invoices, shipping notices, which specify the required data fields for that particular transaction.

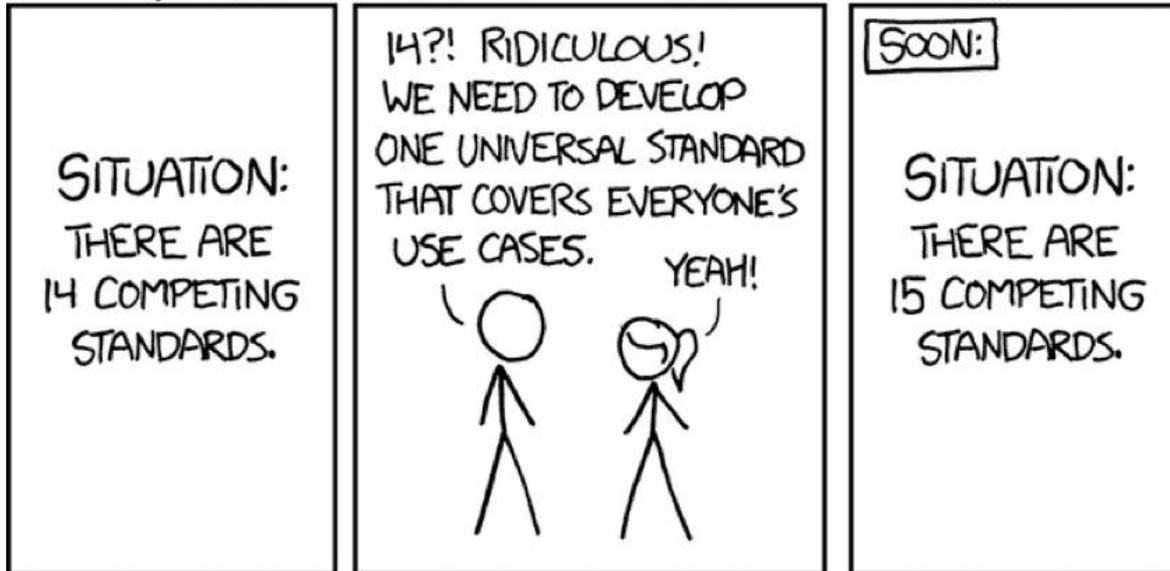
- **Industry-specific subsets:**

In addition to the primary standards, different industries have their own subsets of EDI standards tailored to their needs, e.g.

- EANCOM: retail
- TRADACOMS: Primarily used in the UK retail sector
- ODETTE: automotive industry.
- VDA: European automotive industry, mainly in Germany
- HL7: A semantic interoperability standard commonly used in healthcare data exchange



HOW STANDARDS PROLIFERATE:
(SEE: A/C CHARGERS, CHARACTER ENCODINGS, INSTANT MESSAGING, ETC.)



Ref	Question	A	B	C	D	E
MM1	The role of the Physical data model in the Metadata repository is	Which version of COTS software (E.g. SAP) is implemented	To describe how and where our data is stored in our systems applications or packages.	What the business definition of data concepts is	How many master data records are stored in our MDM system	When the duplicated records were merged
MM2	Updating the Metadata repository is a recommended activity during project close out in the SDLC	TRUE	FALSE			
MM3	What type of Meta-Data provides developers and administrators with knowledge and information about systems?	Unstructured Meta-Data	Process Meta-Data	Business Meta-Data	Data Stewardship Meta-Data	Technical Operational Meta-Data
MM4	These are examples of which type of Meta-Data: Data Stores & Data Involved, Government/ Regulatory Bodies; Roles & Responsibilities; Process Dependencies and Decomposition	Technical Meta-Data	Business Meta-Data	Process Meta-Data	Data Stewardship Meta-Data	Operational Meta-Data
MM5	Which of the following is a Meta-Data scheme focused specifically on documents?	Administrative Meta-Data	Structural Meta-Data	Descriptive Meta-Data	Preservation Meta-Data	Business Meta-Data
MM6	Which of the following contain metrics associated with Meta-Data Management	Steward Representation/ Coverage; Meta-Data Repository Availability; Meta-Data Management Maturity	Meta-Data Management Maturity; No. of Data Stewards; No. of Meta-Data Attributes Listed	Steward Representation/ Coverage; Meta-Data Repository Availability; No. external Reference Data Sources	Meta-Data Repository Availability; No. Data Governance Meetings held annually; No. Meta-Data Repositories	No. of Entities, Attributes & Relationships stored in the Repository

AFTER QUIZ 10

1. General (6)
2. Data Quality (9)
3. Data Storage & Operations (3)
4. Master & Reference Data (9)
5. DW / BI + Big Data (7)
6. Data Modelling (9)
7. Data Security (6)
8. Document, Records & Content Mgt (3)
9. Architecture & lifecycle (4)
10. Metadata management (6)

Maximum possible score = 62

60% (CDMP Associate) = 38

70% (CDMP Practitioner) = 44

80% (CDMP Master) = 50

360

Data Governance

.... at the *Heart* of ALL
Information Management Disciplines

Data Management Process	2%
Big Data	2%
Data Architecture & Lifecycle Management	6%
Document, Records & Content Management	4-5%
Data Ethics	2%
Data Governance	11%
Data Integration & Interoperability	6-7%
Master & Reference Data Management	11%
Data Modelling & Design	11%
Data Quality	11%
Data Security	6%
Data Storage & Operations	4-5%
Data Warehousing & Business Intelligence	11%
Metadata Management	11%



Data Governance and Stewardship (DMBoK 2 revised)

Definition: The exercise of authority, control, and shared decision-making (planning, monitoring, and enforcement) over the management of data assets.

Goals:

1. Enable an organization to manage its data as an asset.
2. Define, approve, communicate, and implement principles, policies, procedures, metrics, tools, and responsibilities for data management.
3. Monitor and guide policy compliance, data usage, and management activities.



Other Definitions of Data Governance

def·i·ni·
The tea
of the r
of an i

“Data Governance is a **quality** control discipline for adding new rigor and **discipline** to the **process** of managing, using, improving and protecting organizational information.”

(IBM Data Governance Council)

“The process of managing and improving data for the benefit of all stakeholders”
(Chris Bradley)

“Data Governance is a **system** of decision rights and **accountabilities** for information-related **processes**, executed according to agreed-upon models which describe who can take what actions with what information, and when, under what circumstances, using what methods.”

(Data Governance Institute)

“Data Governance is the **formal** orchestration of **people, processes, and technology** to enable an organization to leverage data as an enterprise **asset**.”

(MDM Institute)



What Is Data Governance?

The Design & Execution Of Standards & Policies Covering ...

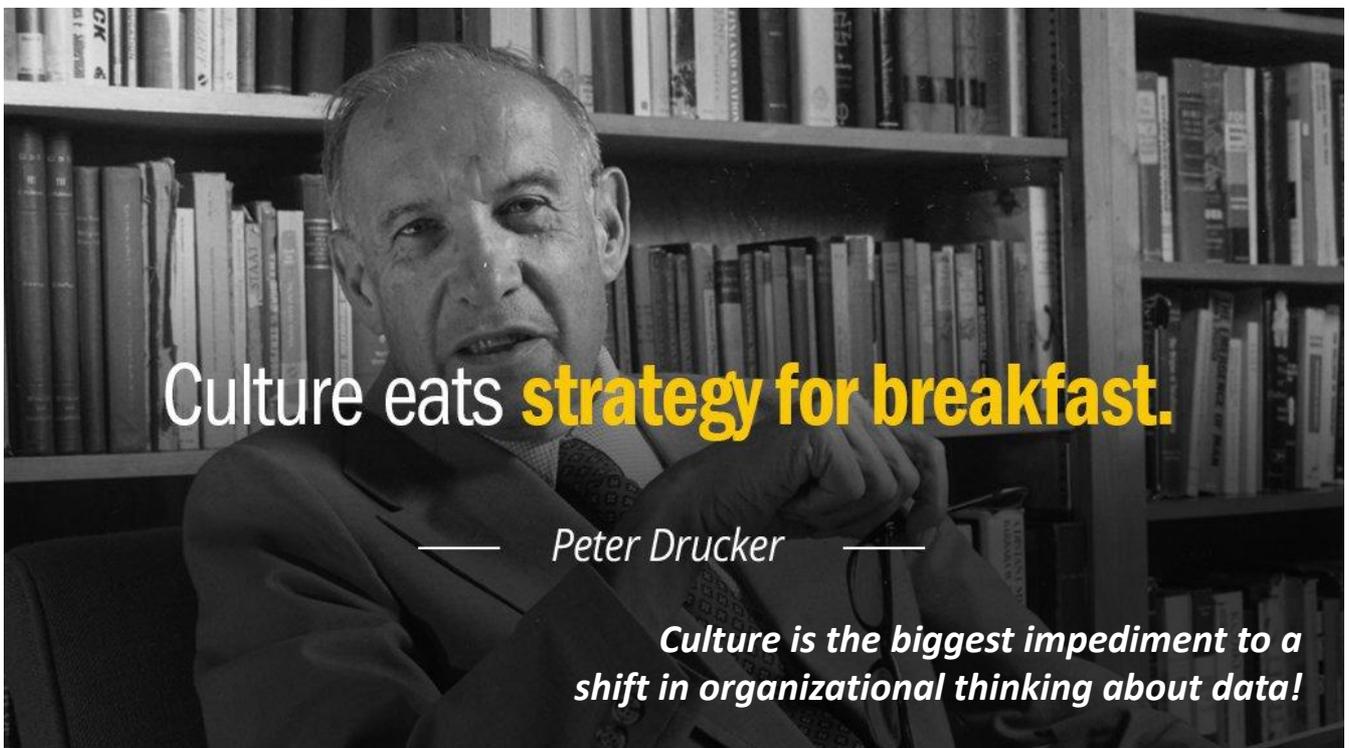
- Design and operation of a management system to assure that data **delivers value**
- Who can do what to the organisation's data and how
- Ensuring standards are set and met
- A strategic & high-level view **across the whole organisation**

To Ensure ...

- Key principles/processes of effective Information Management are put into practice
- Continual improvement via evolution of an Information Management strategy
- Management of Data as a Corporate asset becomes part of the **culture**

Data Governance Is **NOT** ...

- A "one off" Tactical management exercise
- The responsibility of the Technology and IT department alone



Assets

ISO 55000 defines Asset Management as:
the coordinated activity of an organisation to realize value from assets.

ISO 55000 defines Assets as:
An asset is an item, thing or entity that has potential or actual value to an organisation.



Data as an asset

Data is an asset that has value and should be managed and measured accordingly throughout its lifecycle



DATA is our ONLY

- NON Depletable
- Reusable
- Copyable
- Enterprise wide

Asset



P / 368

Key Points

- Data Governance is core (foundational) function of data management
- Data Stewardship comprises business data stewards with data management professionals (IT)
- A Data Governance program requires organizational structure and activities
- Data Governance is a ongoing program

Data Governance: A Core DM Function

- Data Governance interacts and guides how all other data management functions are performed.
- Requires understanding of business strategies and objectives including regulatory constraints to set effective data strategy for all functions.
- Data Governance has a role in the data management activities of all business projects in an enterprise.



P / 369

Data Stewardship

- Comprised of business data stewards and subject matter experts (SMEs) partnered with Data Management professionals (IT) that are accountable for the effective management and use of organization's enterprise data assets.
- *Typical* data stewards' types include:
 - **Executive** – Senior managers serve on Data Governance Steering Committee & DGC (A "Chief Data Steward" or CDO if you have one, typically chairs the Data Governance Steering Committee (DGSC).
 - **Coordinating** – Leads team of business data stewards.
 - **Business** – Recognized data subject matter experts in organization.
- Data stewards' responsibilities can span across all Data Management functions.



P / 370

Type of Data Stewards (DMBoK 2_R)

- **Executive Data Stewards** are senior managers who serve on a Data Governance Council.
- **Enterprise Data Stewards** have oversight of a data domain across business functions.
- **Technical Data Stewards** are IT professionals operating within one of the Knowledge Areas, such as Data Integration Specialists, Database Administrators, Business Intelligence Specialists, Data Quality Analysts, or Metadata Administrators.
- **Chief Data Stewards** may chair Data Governance bodies in lieu of the CDO or may act as a CDO in a virtual (committee-based) or distributed Data Governance organization. They may also be Executive Sponsors.
- **Coordinating Data Stewards** lead and represent teams of business and technical Data Stewards in discussions across teams and with executive Data Stewards. Coordinating Data Stewards is particularly important in large organizations..
- **Business Data Steward**
Business Data Stewards are business professionals, most often recognized subject matter experts, accountable for a subset of data. They work with stakeholders to define and control data.
- **A Data Owner** is a businessperson who is accountable for decisions about data within their domain.



P / 371



Data Owner



- **Accountable** for the Data Subject area for which they are the “owner”
- Work closely with Data Stewards to establish consistent data quality and business rules
- Define & approve quality, access and security requirements.
- Ensure data for the data subject area is fit for purpose
- Communication & Stakeholder management skills are essential
- Corporate **Accountability** for the Data subject area must **not** simply focus on “their” departmental needs.

P / 372





Data Steward

- Is a **business** role appointed to take responsibility for the quality and use of their organization's data assets.
- Develop and maintain definitions
- Authority for data and metadata definitions
- Create ongoing business rules for data
- Support business analysts in the alignment of functional processes and data requirements
- Assist data quality manager in defining metrics, matching and standardisation rules
- Track and determine the need for additional data elements on projects and opportunities for re-use across projects and business processes
- Is a recognized subject matter expert in the data subject area / business domain that he or she is responsible for
- Works collaboratively across the organization with data stakeholders and others Identifying data problems & issues
- Is an effective communicator
- Works in association with the Data Owner to protect and enhance the data assets under his or her control



P / 374

Data Steward - Role Activities

1. Ensuring data is managed in accordance with IMF and data process exist which others can follow
2. Implementing appropriate controls and data governance mechanisms to ensure data is fit for purpose, based on agreed process/lifecycle
3. Defining and monitoring key data to agreed standards and quality (business and regulatory (e.g. BASEL)) needs
4. Assisting in understanding the business data needs to ensure data is fit for all customers (data mapped to processes)
5. Monitoring & enforcing data policies & practices in their business function
6. Communicating & promoting the value of IMF, to ensure data is right first time through appropriate IMF roles
7. Reviewing and monitoring data quality analysis & audits, including MI/outputs from data quality tools
8. Assist in data quality analysis & improvement and prioritise data for remediation
9. Data Triage - Identify and remediate data issues within own business unit or across the Enterprise
10. Providing first point of contact to resolve IM issues
11. Providing input to Data policies, standards & procedures for business function and across the Enterprise
12. Defining and agreeing of business data quality rules for data quality tool

Assist the Data Owners in the implementation & adherence to:

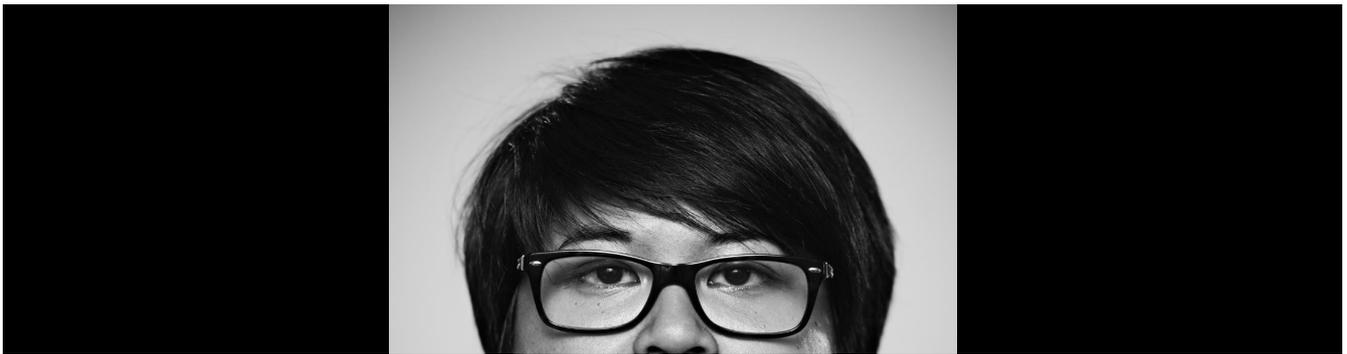
- i. Information Management Principles
- ii. Global Data Privacy Policy
- iii. Global Information and Records Management Procedure
- iv. Group Policy Information Security Management
- v. Group Policy Protection and Disclosure of Price Sensitive Information
- vi. Group Policy Protection and Disclosure of Personal (PII) Information

Other responsibilities

- a) Openly advocating the benefits of a comprehensive Information Management Framework
- b) Educating others on and assisting to implement policies, processes, standards, and procedures
- c) Identifying opportunities to better manage information across functions
- d) Defining, implementing and reviewing appropriate Information Quality standards
- e) Accepting responsibility for the implementation of information quality performance measures
- f) Assisting with gathering and collating feedback for Steering Groups and Data Owners
- g) Collaborating with the wider Information Management community to ensure information quality procedures are correctly designed and followed
- h) Managing the demand for changes to processes, policies, and procedures against budget and resource constraints
- i) Facilitating the prioritisation of work, scheduling and assignment of information quality initiatives
- j) Resolving day-to-day process issues in coordination with other Data Stewards and departmental managers
- k) Assisting with the development training plans, materials and coordinating of training activities related to Information Management



375



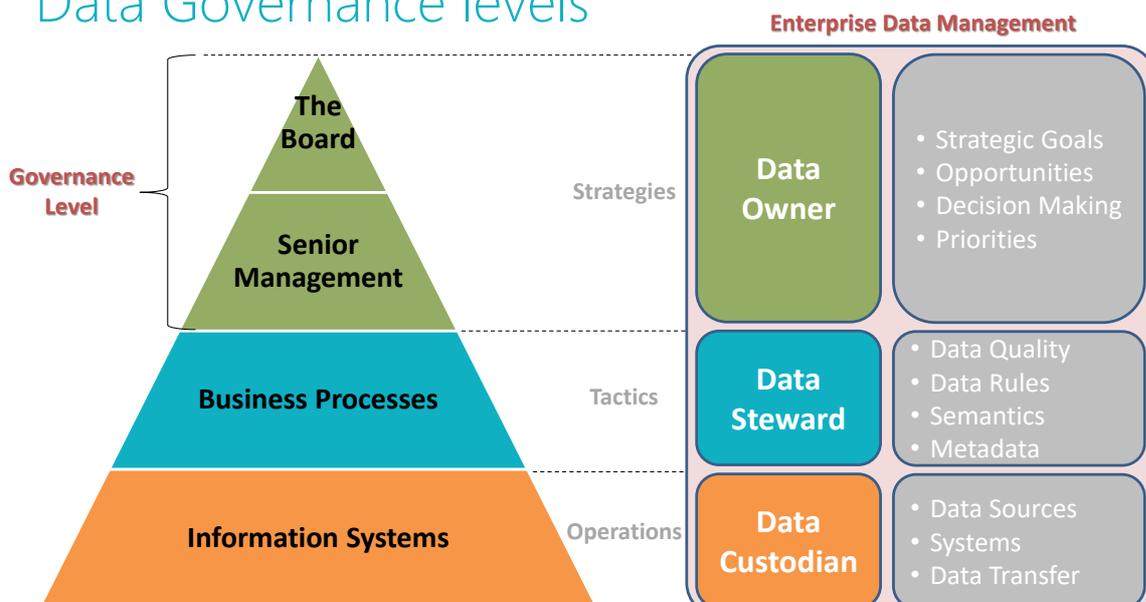
Data Custodian

- Provide formats for operational data and ensure reused
- Assist Data Stewards in interpretation and definition
- Profile and explain source system details
- Work with DQ Manager to correct data at the source
- Explain operational systems processing
- Attend DGC to provide expertise on data sources, targets and transformation processes
- Ensure data is backed up, archived & available to the Business

Container
vs
Content



Data Governance levels

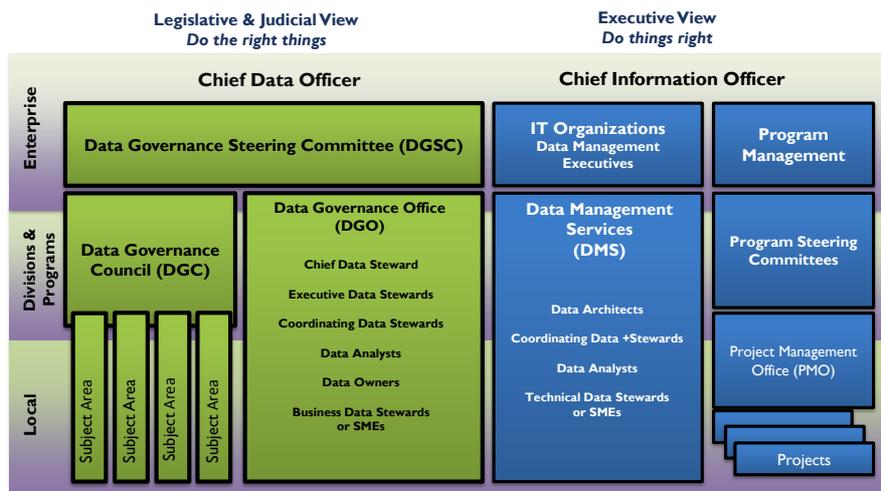


Example Governance (RACI)

DATA SUBJECT AREA <i>(e.g. from CDM)</i>	RESPONSIBLE (R) (DOES)	ACCOUNTABLE (A) (OWNS)	CONSULTED (C)	INFORMED (I)
FINANCE	Geoff Hurst, Data Steward (Finance) Joe Smith, Data Steward (Finance) Joe Cole, Data Steward (Finance) Nobby Stiles, Data Custodian (Finance)	Gordon Banks, Data Owner (Finance)	Alan Ball, Lead Steward (Sales) Alan Knott, Data Owner (Sales)	Bobby Moore, Steward (Marketing)
SALES	Alan Ball, Lead Steward (Sales) Alan Border, Data Steward (Sales) Alan Hunt, Data Steward (Sales)	Alan Knott, Data Owner (Sales)	Gordon Banks, Data Owner (Finance) Alf Ramsey, Data Owner (Marketing)	Bobby Moore, Steward (Marketing)
MARKETING	Bobby Moore, Data Steward (Marketing) David Beckham, Data Steward (Marketing) Ben Stokes, Data Steward (Marketing)	Alf Ramsey, Data Owner (Marketing)	Geoff Hurst, Data Steward (Finance) Alan Ball, Lead Steward (Sales)	Ussain Bolt, Steward (Logistics)
ETC				



Data Management Organizations (DMBoK)



Data Governance includes:

- legislative-like functions (defining policies, standards, and the Enterprise Data Architecture),
- judicial-like functions (issue management and escalation), and
- executive functions (protecting and serving, administrative responsibilities).

To manage risk, organizations need to adopt a representative form of data governance, so that all stakeholders can be heard.

In creating Data Management **Services**, IT involves the Data Governance Council to:

- estimate the enterprise needs for these services,
- provide the justification for staffing and,
- provide the justification for funding to provide these services



Data Management Organizations (DMBoK)

Data Governance Steering Committee	Data Governance Council (DGC)	Local Data Governance Committee(s)	Data Stewardship Teams	Data Governance Office (DGO)
The primary and highest authority organization for data governance in an organization.	Manages data governance initiatives.	Large organizations may have divisional or departmental data governance councils.	Communities of interest focused on one or more specific subject-areas or projects.	Ongoing focus on enterprise-level data definitions and data management standards.
Responsible for oversight, support, and funding of data governance activities. Facilitators responsible for council participation, communication, etc.	Manages development of policies or metrics, issues, and escalations.	Local DGC's work under the backing of an Enterprise DGC.	One or more temporary or permanent focused groups of business data stewards collaborating and consulting with project teams on data definitions and data management standards related to the area of focus.	A staff organization in larger enterprises supporting the efforts of the Data Governance Council, Data Stewardship Steering Committees, and Data Stewardship Teams
Consists of a cross-functional group of senior executives.	Consists of executives according to the operating model used.	Not recommended for smaller organizations.	Consists of business and technical data stewards and data analysts.	Consists of coordinating roles entitled data stewards or custodians, and data owners.
Typically releases funding for data governance and data governance-sponsored activities as recommended by the DGC and CDO.			Led by a coordinating data steward.	Analogous to the PMO (but for data).
This committee may have oversight from higher-level funding or initiative-based steering committees.			Typically, within an assigned subject area	DGO is across ALL Knowledge Areas. i.e. not only DG.



Data Governance Office

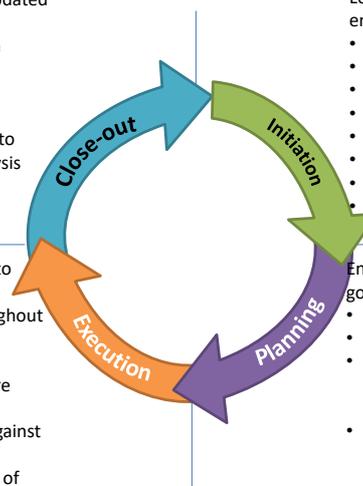
DATA GOVERNANCE HAS TOUCH POINTS THROUGHOUT THE PROJECT LIFECYCLE VIA LIAISON WITH THE DATA GOVERNANCE OFFICE (DGO, SIMILAR TO PMO)

Update Project artifacts at project close out. Updated artifacts should include:

- Updated metadata repository (Business Data Glossary & Technical Data Dictionary)
- Metrics Library
- Data / System Owners and Stewards
- Policies and Guidelines for when users need to take data offline for custom reporting / analysis
- Regulatory reporting
- "Lessons learned" report

Throughout project leverage Data Governance to ensure:

- Appropriate stakeholder participation throughout the project
- Proper vetting of requirements and design documents to ensure data integrity issues are addressed
- Continuous analysis of new metrics / data against existing metrics & data
- Consideration of other business areas usage of this data concept



Leverage Data Governance at project initiation to ensure consistent understanding of:

- Data needs
- Authoritative sources (using Business Glossary)
- Functional usage of data
- Frequency of need
- Needs / use beyond requester
- Stakeholders matrix
- Quality criteria
- Regulatory reporting

Ensure adequate resources / planning for data governance related tasks; e.g.

- Profile to assess integrity of data
- Testing using data quality scenarios
- Training and Communications plan which indicate significance of using the firm's direction for consistent information
- Leverage existing data models



Everyone in the Data Management Community is responsible for communicating and promoting awareness on the value of Data Governance across the organization

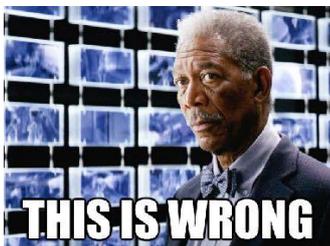


Guiding principles



- Data management is a shared responsibility
- Data Stewards have responsibilities in all 10 (DMBOK 1) management functions (11 in DMBOK 2)
- Every data governance / data stewardship programme is unique
- The best data stewards are found not made
- Shared decision making is the hallmark of data governance
- DG Councils/Data Stewards (legislative) while DMSG (executive)
- Data Governance occurs at enterprise and local levels
- No substitute for visionary and active IT leadership
- Centralised organisation for DM professionals is essential
- Define a formal charter for the Data Governance Council
- Data Strategy should be driven by the Business Strategy

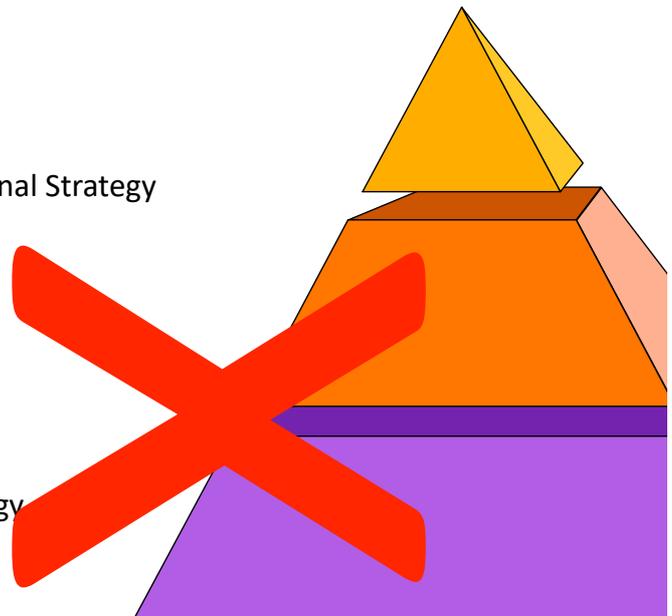
Data Strategy in Context



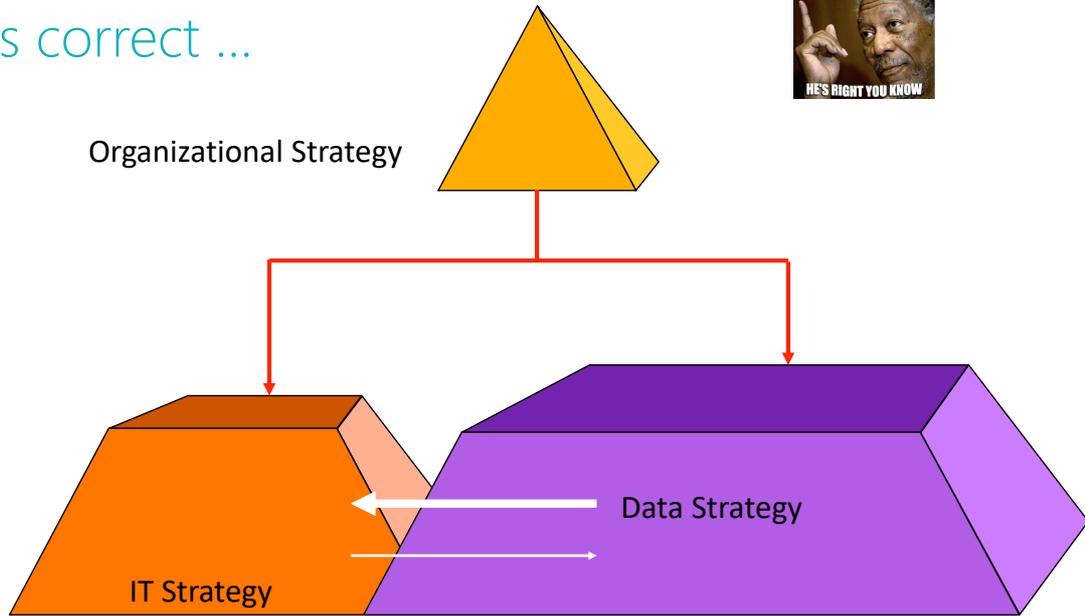
Organizational Strategy

IT Strategy

Data Strategy



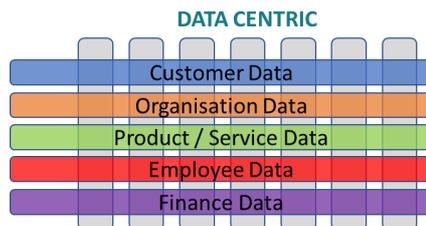
This is correct ...



Organisational Models For DG

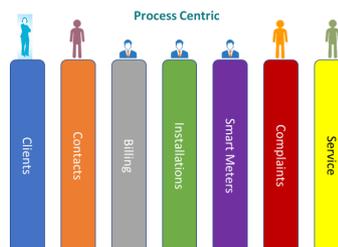
DATA CENTRIC

Business appointed full time or part time roles accountable for improvement of key data domains wherever created or used across an organisation, e.g., Data Stewardship.

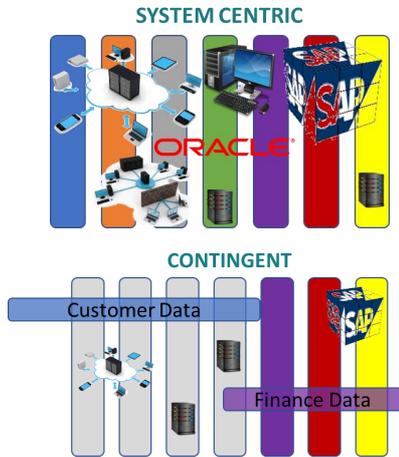


PROCESS CENTRIC

Process owners become the data owner for all data created, amended & deleted by the business process for which they are responsible.



Organisational Models For DG



SYSTEMS CENTRIC

System owners become the data owner for all data created, amended & deleted by the system for which they are responsible.

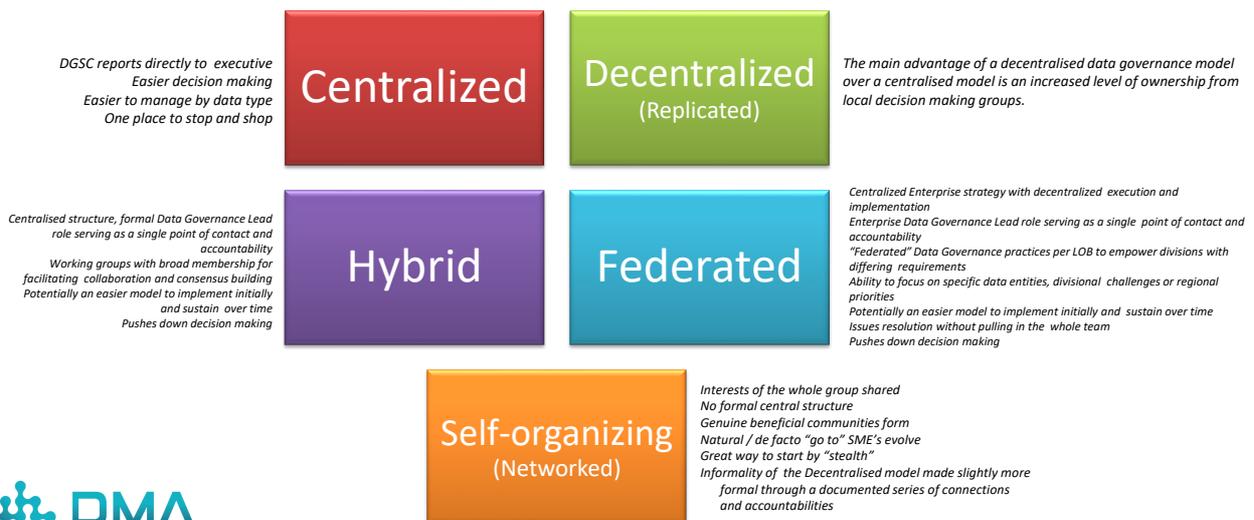
CONTINGENT / HYBRID

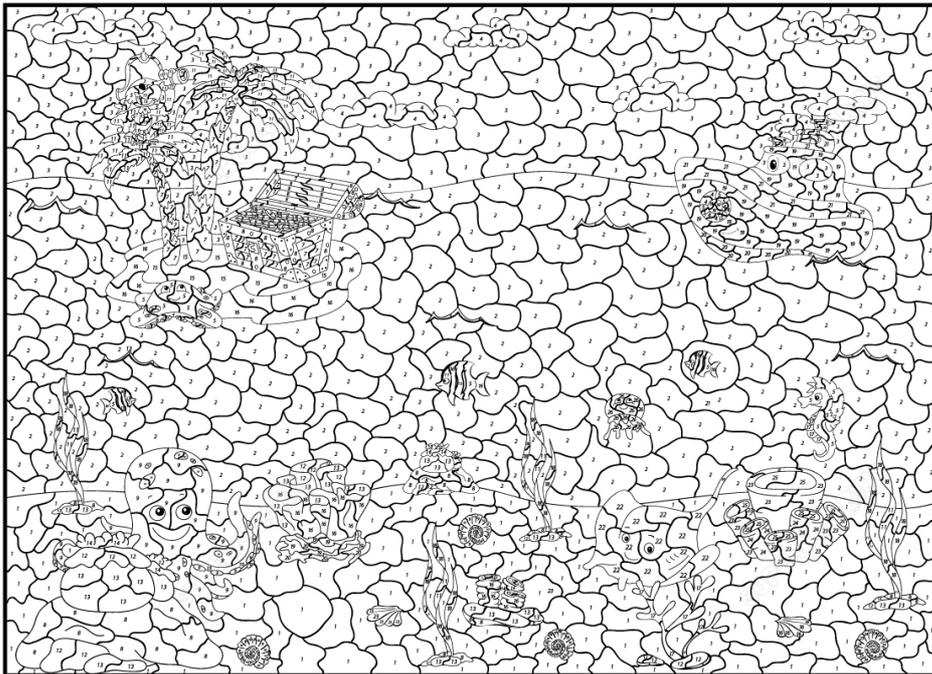
There is no single best model for data governance, either when initiating data improvement activities, or as Business As Usual.

The 'best' model is dependent on the type of data and the circumstances of each initiative, at each stage of maturity



Operating Models For DG





- 1
- 2
- 3
- 4
- 5
- 6
- 7
- 8
- 9
- 10
- 11
- 12
- 13
- 14
- 15
- 16
- 17
- 18
- 19
- 20
- 21
- 22
- 23
- 24
- 25
- 26
- 27



print the world

Start Small – DG will ripple out



**“Data Definitions for our
KYC(CDE) items”**



**What are the top
30 CDE data items
to be defined?**

→ **High level
model(s)**

**Who else needs to be
Consulted, Informed
about these items?**

→ **RACI**

**What are the current data
problems / Who decides if the
definitions are conflicting?**

→ **Ownership /
Use cases**

**What are the
quality criteria for
these items?**

→ **DQ /
Ownership**

**How do Data
Definition owners
undertake their role?**

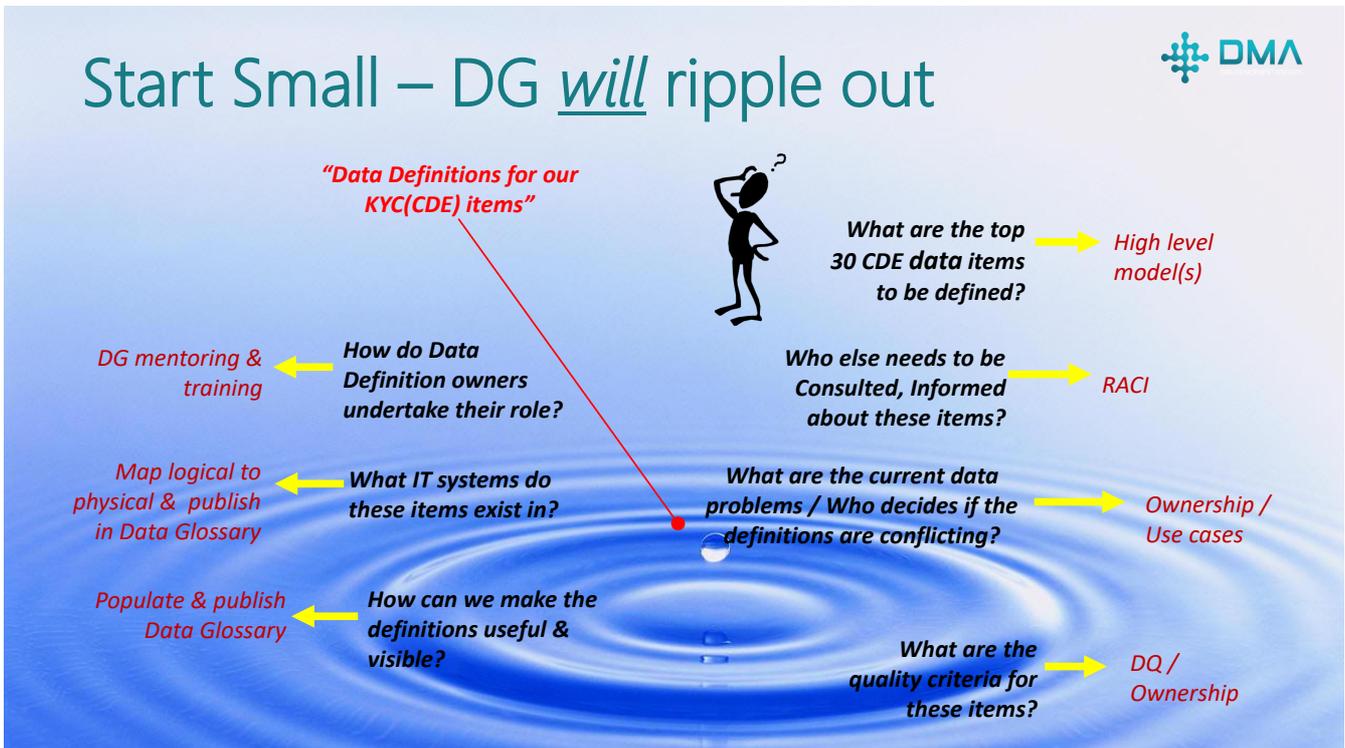
← **DG mentoring &
training**

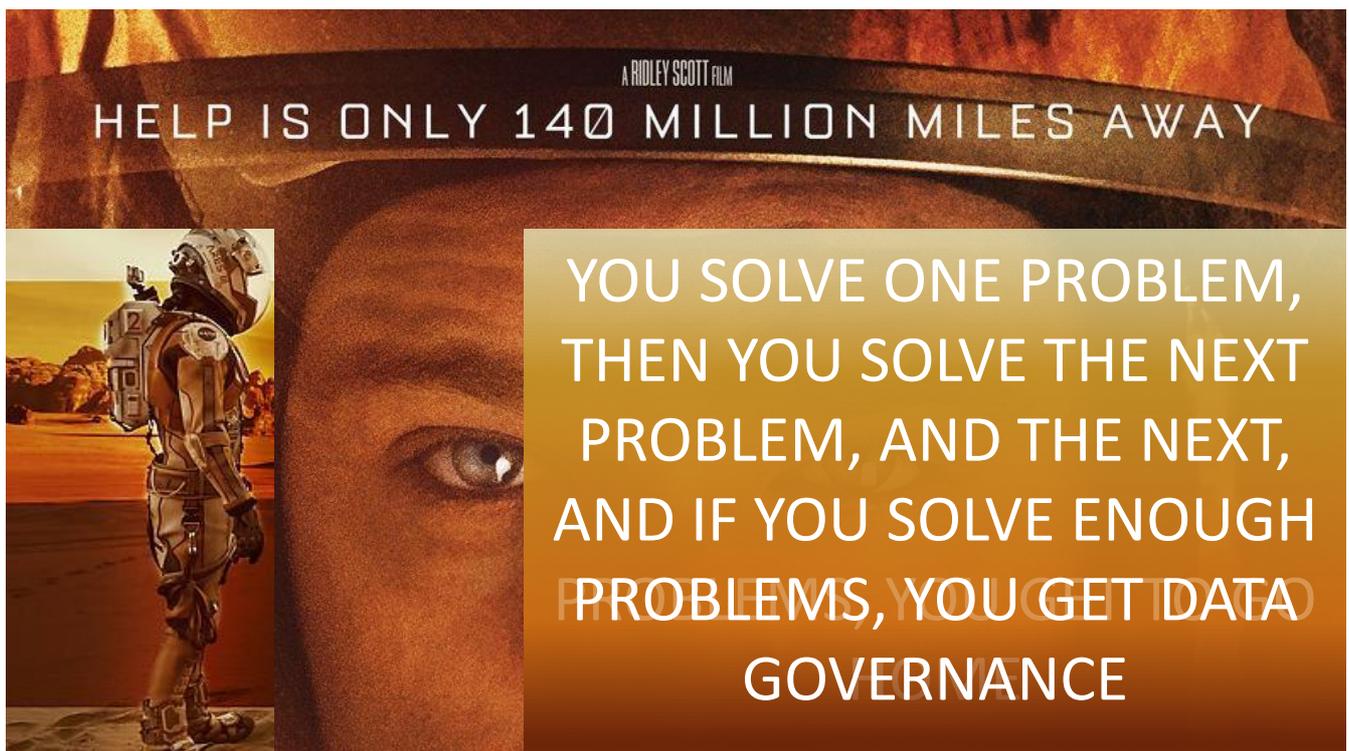
**What IT systems do
these items exist in?**

← **Map logical to
physical & publish
in Data Glossary**

**How can we make the
definitions useful &
visible?**

← **Populate & publish
Data Glossary**





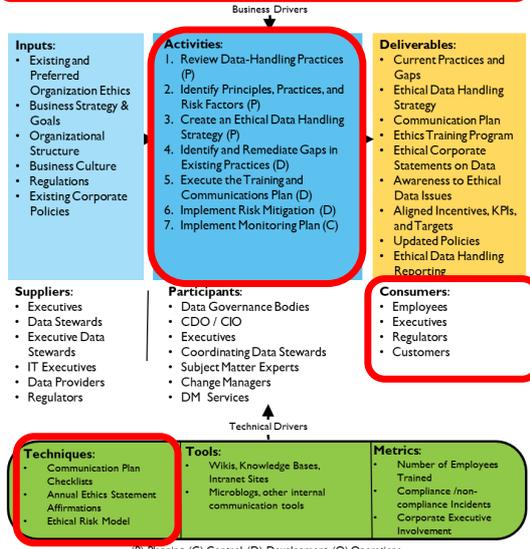
Data Governance Summary

- Data Governance is a **continuous process** of data improvement
- Data Governance is a **Business change programme**
- **IT** is a key stakeholder in DG, but **is NOT responsible for DG**
- There are **different organization models** for DG from Centralized, through Hybrid, to Federated & Self Organising Teams &
- DG is the exercise of authority and control over the management of **data assets**
- The 3 most important CSF's for Data Governance success are **Communication, Communication, Communication**
- **Everyone** in the Data Management Community is responsible for communicating and promoting awareness on the value of Data Governance across the organization
- The best **Data Stewards are found**, not made

Data Handling Ethics (DMBoK 2 Revised)

Definition: Data handling ethics are concerned with how to procure, store, manage, interpret, analyze / apply and dispose of data in ways that are aligned with ethical principles, including potential impacts on people and organizations.

Goals:
 1. To control risk by establishing acceptable practices for data handling.
 2. To change/instill preferred culture and behaviors on handling data.
 3. To maintain alignment with compliant practices in ethical handling of data.



Data Ethics

Data Handling Ethics are based on medical / bio ethics*

Respect for

1. autonomy,
2. beneficence,
3. non-maleficence, and
4. justice

... referred to as the four pillars of medical ethics

The DMBoK2 has 3 !

The Beneficence Principle has two elements:

1. Do no harm
2. Maximize possible benefits and minimize possible harms

The Justice Principle considers the fair and equitable treatment of people.

The Respect for Persons Principle reflects the fundamental ethical requirement that people be treated in a way that respects their dignity and autonomy as human individuals.

* The Belmont Principles were developed by the U.S. HHS in 1979 to guide the ethics of medical research, and they are also applicable as guiding principles within the field of Data Management

European Commission Article 29 Data Protection Working Party criteria to evaluate anonymization methods.

*Anonymity is protected when it is **only** possible to analyse sizeable "clusters" of individuals who **cannot** be distinguished from one another based on their attributes.*



Data Ethics

Data Handling Ethics are seeking to:

- avoid loss of reputation for the organization &
- loss of customers.

Misleading Visualisations

e.g. Leaving data points out, comparing two facts without clarifying their relationship, or ignoring accepted visual conventions, changing scale to make a trend line look better or worse.



The 5 “C’s” of Data handling ethics ..

- Consent,
- Clarity,
- Consistency,
- Control (and transparency),
- Consequences (and harm)

Data Ethics

- Core **concepts** of data handling ethics are:
 1. Impact on people,
 2. the potential for misuse, &
 3. the economic value of data
- Core **deliverables** for introducing an **ethical data handling capability**:
 1. Current Practices and Gaps,
 2. Ethical Data Handling Strategy ,
 3. Communication Plan ,
 4. Ethics Training Program ,
 5. Ethical Corporate Statements on Data
- **Social licence** is the alignment between:
 1. stakeholder expectations and
 2. the organisation
- **Consent** is considered to be given when an organisation has: obtained explicit permission from an individual for the collection, use or disclosure of personal information.



- **Principles** to understand and evaluate **ethical use of data** within an organization should be based upon:
 1. fairness,
 2. respect,
 3. responsibility,
 4. integrity,
 5. quality,
 6. reliability,
 7. transparency and
 8. trust
- **Ethical responsibility** for managing data to:
 1. reduce risks of misrepresentation,
 2. misuse or
 3. misunderstandinglies with Data Management professionals to:
 4. manage data and
 5. manage the associated risks



Ref	Question	A	B	C	D	E
DG1	According to the DAMA DMBOK, the Data Governance steering committee (aka Data Governance Council (DGC)) is the primary and highest authority organization for data governance in an organization. Who should typically chair this Council?	The Chief Information Officer (CIO)	Chief Data Steward (Business) / Chief Data Officer	The chair should rotate across the Data Owners	The Chief Data Architect	Any Executive / C-level participant in the DGC
DG2	What are the primary characteristics of a data steward?	A business role appointed to take responsibility for the quality and use of their organization's data assets.	Analyzing data quality	The manager responsible for writing policies and standards that define the data management program for an organization.	Identifying data problems & issues	The data analyst who is the subject matter expert (SME) on a set of reference data.
DG3	Which of these is NOT true of Data Governance?	DG is a continuous process of data improvement	IT is a key stakeholder in DG	A DG initiative should always be led by the IT department	There are different organization models for DG	DG is the exercise of authority and control over the management of data assets
DG4	Who is responsible for communicating and promoting awareness on the value of Data Governance in the organisation?	Central Communications and Corporate Awareness	Data Stewards	The Chief Executive Officer	Senior Management Executive Forum	Everyone in the Data Management Community
DG5	Communicating the value of Data Governance can be approached in a number of ways. Which of the following approaches is NOT a recognised way of doing this?	Providing only negative communications on ongoing data issues to key executive stakeholders	Maintaining an intranet website	Publishing a regular newsletter via hardcopy or email	Promote participation in a DM forum or community	Creating a series of "elevator pitches" for the appropriate audience
DG6	When considering a Data Governance program, communication is a key element. There are many ways of managing this communication, with one of the most effective being a Data Management intranet. Which of the following would you typically NOT put onto such a communication vehicle	Description of the DG organisation, it's key members and contact details	Executive message regarding significant data management issues	The Data Steward team profiles	Raw data results of an investigation into a possible data privacy breach	Link to a "raise an issue" log
DG7	How can the Data Governance process in an organization best support the requirements of various regulatory reporting needs?	by providing a business glossary based lookup facility for data definitions	by ensuring that data is properly categorized, owned, understood, defined, documented, and controlled	by performing an as-is data audit	by highlighting the challenges of multiple data definitions within the enterprise	by creating a map of where the enterprise data is located in IT systems
DE1	One way of defining ethics is:	"doing it wrong, and then apologising".	"doing it right when someone is looking".	"doing it wrong, and then expertly covering it up".	"doing it wrong, and failing to covering it up".	"doing it right when no one is looking".
DE2	The ethics of data handling, centre on several core concepts. They are:	impact on people, potential for misuse and economic value.	impact on people, potential for re-use and storage cost.	access to data, potential for misuse and storage cost.	accurate business glossary, data quality and reference data.	privacy, security and authorisation.
DE3	In data handling ethics, 'social licence' refers to the alignment between:	stakeholder demands and technology deliverables.	stakeholder expectations and the organisation	public perception and published fact.	society's needs and their right to access data.	social and political decision matrices.

AFTER QUIZ 11

1. General (6)
2. Data Quality (9)
3. Data Storage & Operations (3)
4. Master & Reference Data (9)
5. DW / BI + Big Data (7)
6. Data Modelling (9)
7. Data Security (6)
8. Document, Records & Content Mgt (3)
9. Architecture & lifecycle (4)
10. Metadata management (6)
11. Data Governance & Ethics (10)

Maximum possible score = 72

60% (CDMP Associate) = 44

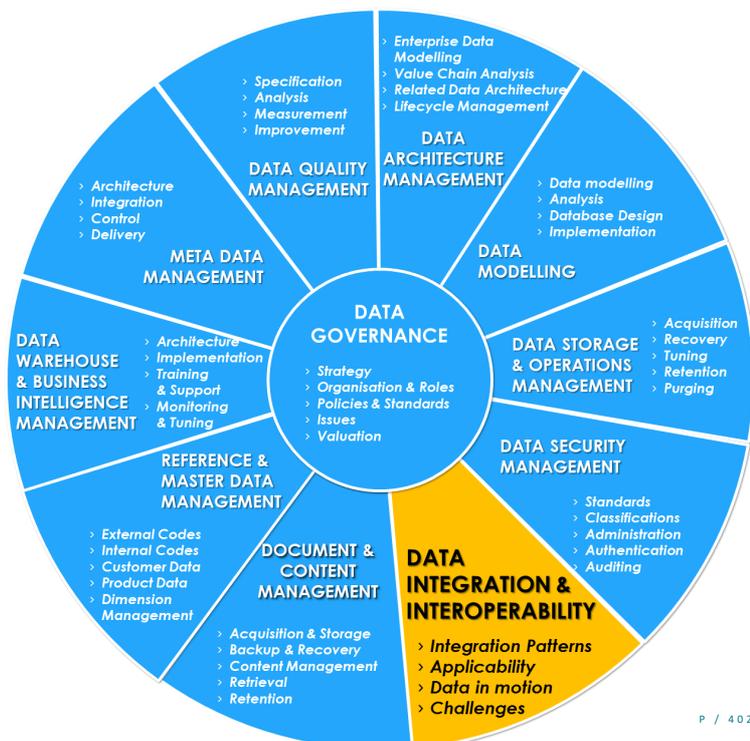
70% (CDMP Practitioner) = 51

80% (CDMP Master) = 58

401

Data Integration & Interoperability

Data Management Process	2%
Big Data	2%
Data Architecture & Lifecycle Management	6%
Document, Records & Content Management	4-5%
Data Ethics	2%
Data Governance	11%
Data Integration & Interoperability	6-7%
Master & Reference Data Management	11%
Data Modelling & Design	11%
Data Quality	11%
Data Security	6%
Data Storage & Operations	4-5%
Data Warehousing & Business Intelligence	11%
Metadata Management	11%



Data Integration and Interoperability (DMBoK 2 revised)

Definition: Data Integration is the movement and consolidation of data within and between data stores, applications and organizations. Data Interoperability is the ability for multiple systems to communicate.

- Goals:**
1. Make data available from disparate sources in the format and timeframe needed.
 2. Lower cost and complexity of managing solutions by developing shared models and interfaces.
 3. Identify meaningful events and automatically trigger alerts and actions.
 4. Support business intelligence, analytics, master data management, and operational efficiency efforts.



Goals

Applying Data I&I practices and solutions aims to:

- Make data available in the format and timeframe needed by data consumers, both human and system
- Consolidate data physically and virtually into data hubs
- Lower cost and complexity of managing solutions by developing shared models and interfaces
- Identify meaningful events (opportunities and threats) and automatically trigger alerts and actions
- Support business intelligence, analytics, master data management, and operational efficiency efforts



Core Concepts:

Integration: Movement and consolidation of data into consistent forms (physical or virtual) within and between data stores, applications, and organizations.

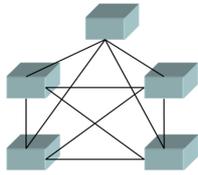
Interoperability: The ability for multiple systems to communicate through data that does not require processing on either side

Common use cases:

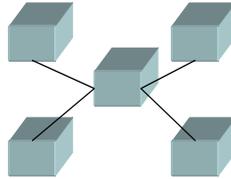
1. Integration of data between data stores, including ingesting from and disbursing data to external sources
2. Consolidation of data stores, including application consolidation, data hub management, and mergers and acquisitions
3. Distribution across data stores and data centres
4. Moving data into archives, and updating data from one archive technology to another to enable future use



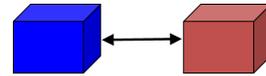
Point to Point



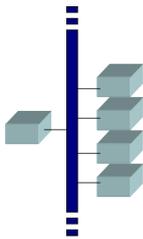
Hub Distribution



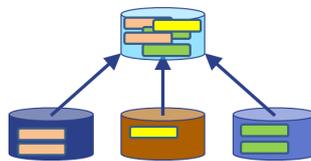
Message Synchronisation & Propagation



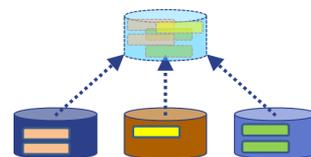
Bus Distribution



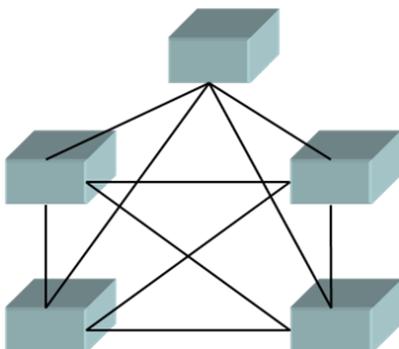
ETL, ELT & CDC



Abstraction / Virtual Consolidation



Point to Point



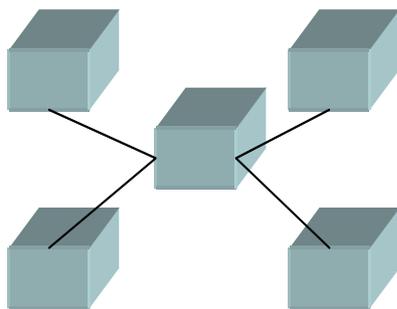
- Good for a very limited number of interfaces
- Short life expectancy of interfaces
- Need for speed

Issues:

- Approach does not easily / predictably scale.
- Run-time: performance issues may arise when many systems want to get the same data from the same source.
- Design-time: issues arise when multiple systems require different versions or formats of the data.
 - The source systems support teams must build or support separate interfaces for each specific use.
- Support: While point-to-point interfaces may be most performant, the whole ecosystem may suffer because exterior effects are not considered.



Hub Distribution

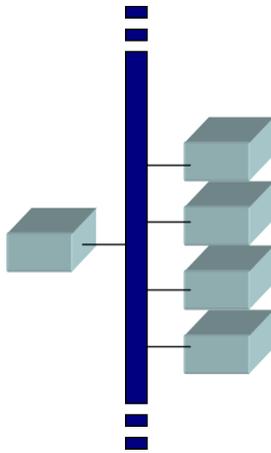


- Data updates (published) into a central hub repository which distribute them to all authorized (subscribing) applications and may also keep a copy in a repository.
- An architecture formerly used by Master & Reference Data Management systems (MDM); if the messaging hub is modified to persist data, the hub can act as the Record of Reference.
- Messaging hub / repository interprets each transaction, knows how to transform it, and which applications are authorized to receive the updates.

Potential Issues:

- Situations (e.g., security or regulations) where data should not be shared in a hub. Use a small dedicated secured hub for that data.
- Latency may be unacceptable, or performance is insufficient.
- Single point of failure

Bus Distribution

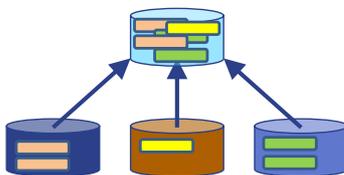


- Bus Distribution model “pushes” data into a central service which then pushes it to authorized applications or the application can “pull” the data.
- Data is not permanently retained and therefore it cannot act as the Record of Reference.
- Often described as publish-subscribe.
- Basis of a Services Oriented Architecture (SOA).
- Well-suited for applications using a one-to-many distribution which do not have data retention requirements for the distribution tool
- Compared to the Hub model, this model is more scalable, can offer better performance, and is well-suited for a multi-mega-centre environment.
- It has higher initial build costs but with lower maintenance costs.
- Same single point of failure concern as the Hub model.

P / 411

Physical movement & consolidation

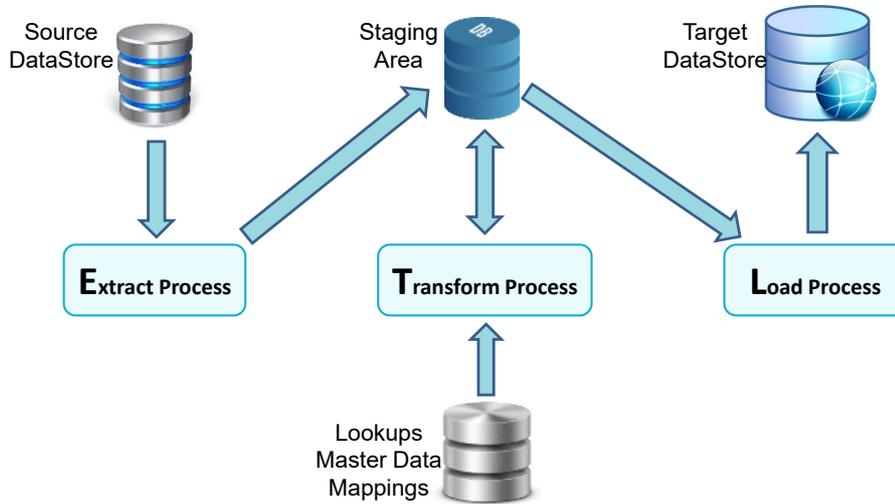
ETL, ELT & CDC



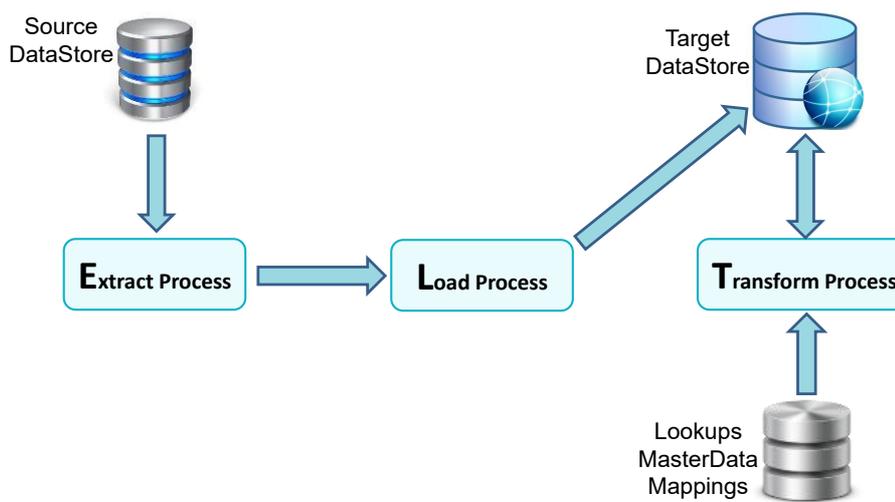
- Batch distribution for the mass movement of data collected over time from a source data structure and distributed as a batch to another.
- Extract Transform & Load (ETL) and Extract Load and Transform (ELT tools facilitate this process with scheduling, parallel processing and complex data transformation, cross reference and data mapping capabilities.
- Batch may contain incremental or full downloads of the data. This is probably the most commonly used approach to Data Integration.
- ETL / ELT is typically run according to a schedule (e.g. the “overnight load of a Data Warehouse) and is used for bulk data movement, usually in in batch.
- Change Data Capture (CDC) is event driven (e.g. the stock level has fallen below critical) and delivers real-time incremental replication.

P / 412

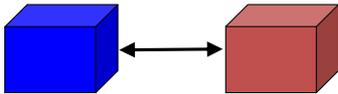
ETL Process Flow



ELT Process Flow



Message Synchronisation & Propagation



- Message based synchronisation and data propagation is used for application to application integration.

2 main genres:

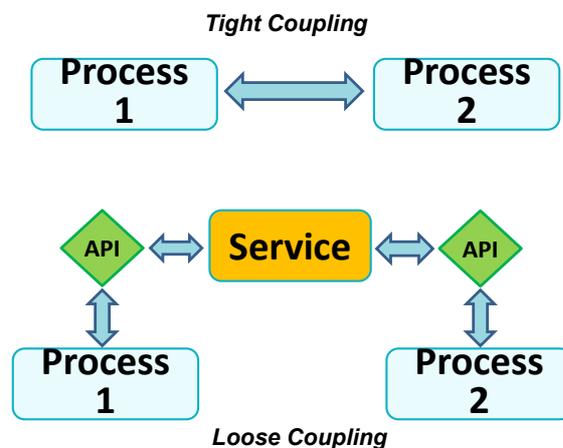
- Enterprise Application Integration (EAI)
- Enterprise Service Bus (ESB)
- Can be based on a Hub or a Bus.
- Both of these are used primarily for the purpose of event driven business process automation.
- May be loosely or tightly coupled

Application Coupling

Coupling describes how an interface connects any two systems.

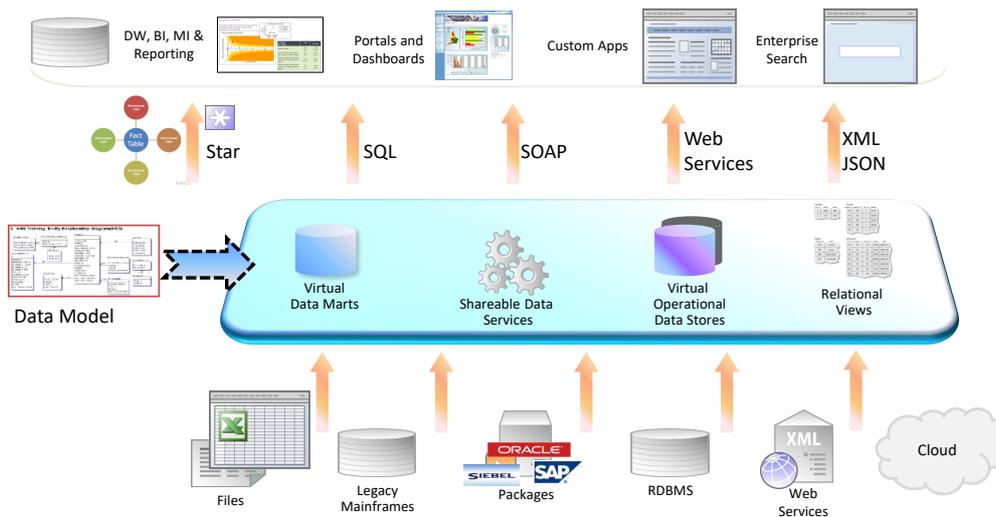
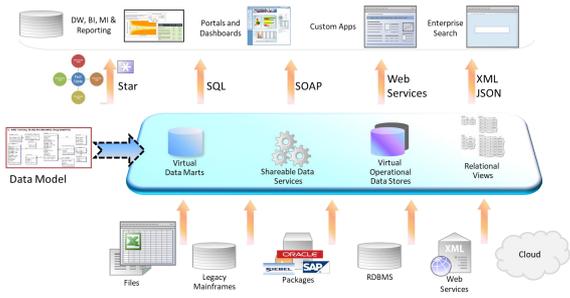
Tight coupling means that the sources and targets are pre-defined and directly named.

Loose coupling allows the sources and targets to **not be linked directly**; instead, a service is used to allow sources and targets to remain anonymous to each other, and instead, the applications use APIs to link to the service.



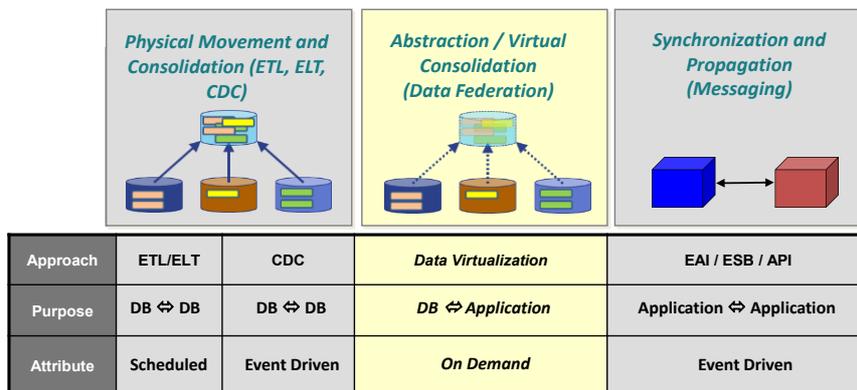
Abstraction / Virtual Consolidation

- Virtual distribution leaves data in situ while building an intelligent layer on top of it that can locate and transform data on the fly.
- Often referred to as Enterprise Information Integration (EII)
- No need to touch underlying source data
- Very rapid time to solution
- Easily mix structured & unstructured data sources
- Present data to applications as SOAP, SQL, STAR, ODBC, XML, JSON
- Rapidly emerging option for organizations with large legacies that they intend to retain without making major investments in transforming them.
- Consuming applications see data as though it's in their own systems





Key Differences with Data Virtualisation



Other DII Terms & Considerations



Latency: The **time difference** between when **data is generated** in the source system and when the **data is available** for use in the target system. Different approaches to data processing result in different degrees of data latency. Latency can be high (batch) or low (event-driven) to very low (real-time synchronous).



Asynchronous Data Flow: This allows the source system to **continue processing without waiting** for the receiving system's acknowledgment. It implies that **systems can operate independently**, and updates are near-real-time, usually seconds or minutes apart. This approach does not disrupt the source application's processing if target applications are unavailable. However, it results in a delay between updates made in the source and those relayed to target data sets.



Synchronous integration: This requires an **executing process to wait for confirmation** from other applications before proceeding. It **ensures data synchronization** but can process fewer transactions due to waiting times. If any application required for the update is unavailable, the transaction cannot be completed. This method keeps data sets perfectly in sync, often using database capabilities like two-phase commits, especially in financial institutions to synchronize transaction and balance tables.

Ref	Question	A	B	C	D	E
DII 1	The need to manage data movement efficiently is a primary driver for:	Data Integration and Interoperability	Data Storage and Operations	Data Warehousing and Business Intelligence	Document and Content Management	Data Security
DII 2	The acronym ETL most commonly stands for:	Export Transform Log	Extract Transform Load	Extend Trim Load	Efficient Trace Logging	Extract Transpose Leverage
DII 3	Mapping requirements and rules for moving data from source to target enables:	load	extract	transformation	analysis	backups
DII 4	When integrating two data stores using batch or real-time synchronous approaches, results in a difference in:	data quality	lethargy	source of truth	timestamping	latency
DII 5	If two data stores are able to be inconsistent during normal operations, then the integration approach is:	Streaming	Synchronous	Faulty	Asynchronous	Uncontrolled
DII 6	A 'Content Distribution Network' supporting a multi-national website is likely to use:	a replication solution	an extract transform and load solution	a database backup and restore solution	an archiving solution	a records disposal solution
DII 7	Data that is used infrequently or not at all may be moved to an alternative data store. This is called:	replication	analysis	archiving	auditing	authentication
DII 8	Three common interaction models for data integration are:	point to point, hub and spoke, publish and subscribe.	point to point, wheel and spoke, public and share.	plane to point, harvest and seed, publish and subscribe.	straight copy, curved copy, roundabout copy.	record and pass, copy and send, read and write.

AFTER QUIZ 12

1. General (6)
2. Data Quality (9)
3. Data Storage & Operations (3)
4. Master & Reference Data (9)
5. DW / BI + Big Data (7)
6. Data Modelling (9)
7. Data Security (6)
8. Document, Records & Content Mgt (3)
9. Architecture & lifecycle (4)
10. Metadata management (6)
11. Data Governance & Ethics (10)
12. Data Integration & Interoperability (8)



Maximum possible score = 80

60% (CDMP Associate) = 48

70% (CDMP Practitioner) = 56

80% (CDMP Master) = 64

P / 424

DMBoK₂ and CDMP® Preparation Live and On-Demand Classes



Data Management Fundamentals

<https://www.dataversity.net/dmbok-and-cdmp-preparation-learning-plan/>

Data Governance

<https://training.dataversity.net/learning-paths/dgs0-dmbok-and-cdmp-preparation-data-governance-specialist-learning-plan>

Data Modelling

<https://training.dataversity.net/learning-paths/dms0-dmbok-and-cdmp-preparation-data-modeling-specialist-learning-plan>

Data Quality Management

<https://training.dataversity.net/learning-paths/dqs0-dmbok-and-cdmp-preparation-data-quality-specialist-learning-plan>

Master & Reference Data Management

<https://www.dataversity.net/>

Metadata Management

<https://www.dataversity.net/>

 CDMP Online learning program

 Approved by DAMA-I

 Part of DATAVERSITY training center

 Based on DMBoK CDMP syllabus

20% discount code for
DAMA UK
"dmadvisors"

P / 425



Contact

-  info@DMAdvisors.co.uk
-  +44 7973 184475 (mobile) +44 1225 923000 (office)
-  [@inforacer](https://twitter.com/inforacer)
-  uk.linkedin.com/in/christophermichaelbradley/
-  infomanagementlifeandpetrol.blogspot.com



TRAINING
ADVISORY
CONSULTING
CERTIFICATION

