

# Data Cheat Sheet

Term	Concept	Description /Analogy	Illustration	Real world Example	Relevance to Ki
<b>Systems of Record</b>	An information storage system that is the authoritative data source for a given data element or piece of information	<ul style="list-style-type: none"> <li>Used as a data repository where the data object, as a whole or specific attributes of a data object, are maintained (data creation, updating, modifying, and deleting)</li> <li>Helps organisations to manage big data by serving as a primary source</li> </ul>		<a href="#">Electronic Health Records</a>	
<b>Systems of Insight</b>	A system of insight facilitates organizing, transforming, and analysing your data through the consumption, collection, and analysis of data from the combined sources of traditional “systems of record” and “systems of engagement”	<ul style="list-style-type: none"> <li>By leveraging behaviour-driven insights, systems of insight apply advanced analytics to generate customer interaction data</li> <li>Allows businesses to immediately deploy personalisation</li> </ul>		<a href="#">Analyse business data and social feeds, such as social networks, customer service interactions, web clickstream data to predict potential customer behaviour and determine targeted ads</a>	
<b>Systems of Engagement</b>	A system that manages and promotes user collaboration and interaction.	<ul style="list-style-type: none"> <li>A system of engagement overlays and complements an organisation’s investment in a system of record by providing easy access to data, as well as easy to use applications that enable collaboration across your organisation</li> </ul>		<a href="#">Enhance patient engagement through mobile notifications system</a>	
<b>Data Warehouse</b>	A central repository of information that can be analysed to make more informed decisions	<ul style="list-style-type: none"> <li>Data flows into a data warehouse from transactional systems, relational databases, and other sources, typically on a regular cadence</li> <li>Enables powerful analytics on huge volumes (petabytes and petabytes) of historical data in ways that a standard database cannot</li> </ul>		<a href="#">Google BigQuery</a>  <a href="#">SQL Data Warehouse</a>	
<b>Data Lake</b>	A storage repository that holds a vast amount of raw data in its native format until it is needed	<ul style="list-style-type: none"> <li>Enables storage of all structured or unstructured data at scale</li> <li>Uses a flat architecture to store data</li> <li>Ability to run different types of analytics to guide better decisions</li> </ul>		<a href="#">Azure Data Lake</a>	

Term	Concept	Description /Analogy	Illustration	Real world Example	Relevance to Ki
<b>Data Mesh</b>	Data mesh is an architectural paradigm that embraces the ubiquity of data in the enterprise by leveraging a domain-oriented, self-serve design	<ul style="list-style-type: none"> <li>Used to stitch together data held across multiple data silos</li> <li>Allows for a federated data architecture by bringing autonomy to data domains, while providing standardisation through central capabilities</li> </ul>		<a href="#">JPMorgan Chase build a data mesh architecture to enable data sharing across the enterprise while giving data owners the control and visibility they need to manage their data effectively</a>	
<b>Centralised data platform approach</b>	Data brought into a central platform to enable insights and analytics to allow users to use the data for strategic purposes	<ul style="list-style-type: none"> <li>Enables a cohesive view of data from multiple sources</li> <li>Facilities better optimisation through better tracking, improved data integrity, reduced data redundancy, and more</li> </ul>		<a href="#">By taking a centralised data platform approach, the Bank of England was able to break down data siloes, improve speed of reporting by more than 500x from days to minutes, and ensure data is reliable and compliant</a>	
<b>Data streaming</b>	The process of transmitting, ingesting, and processing data continuously rather than in batches	<ul style="list-style-type: none"> <li>Enables generation of analytic results in real time as streamed data can be processed and analysed as it arrived</li> <li>Benefits include easy data scalability, detection of patterns in time-series, and more</li> </ul>		<a href="#">Real-time stock trades, up-to-the-minute retail inventory management, social media feeds, multiplayer game interactions, and ride-sharing apps</a>	
<b>Data model</b>	The process of creating a visual representation of either a whole information system or parts of it to communicate connections between data points and structures	<ul style="list-style-type: none"> <li>Illustrates the types of data used and stored within the system, the relationships among these data types, the ways the data can be grouped and organised and its formats and attributes</li> </ul>		<i>See below data modelling techniques for examples</i>	
<b>Hierarchical Model</b>	A tree-like structure where there is one root node (parent node) and the other child nodes are sorted in a particular order	<ul style="list-style-type: none"> <li>Simple and rigid structure</li> <li>Very rarely used nowadays as retrieving and accessing data is difficult in a hierarchical database</li> </ul>		<a href="#">Real-world relationships e.g. food recipes, sitemap of a website</a>	
<b>Object-oriented Model</b>	A data model that consists of a collection of objects, each with its own features and methods	<ul style="list-style-type: none"> <li>Communicates while supporting data abstraction, inheritance, and encapsulation</li> </ul>		<a href="#">db4o, Smalltalk and Cache</a>	

Term	Concept	Description /Analogy	Illustration	Real world Example	Relevance to Ki
<b>Network Model</b>	A data model in which each record can be linked with multiple parent records	<ul style="list-style-type: none"> <li>Similar to hierarchical model but easier to convey complex relationships</li> <li>Useful for projects of a complex nature or where activities are subject to a considerable degree of uncertainty in performance time</li> </ul>		<a href="#">Modelling protein-protein interaction during the drug discovery process for the pharmaceutical industry</a>	
<b>Entity-relationship Model</b>	A high-level data model that defines data entities (e.g. concept, data, object) and their relationship for a specified software system	<ul style="list-style-type: none"> <li>Represented in a diagram called an entity-relationship diagram, consisting of Entities, Attributes, and Relationships</li> <li>Provides a better view of the data that is easy for stakeholders and developers to understand</li> </ul>		<a href="#">Database of a school: Entities - 'Teacher', 'Department' Attributes - (of 'Teacher') Salary, Age, ID Relationship - 'Teacher' works for 'Department'</a>	
<b>Relational Model</b>	A data model used to describe the different relationships between the entities	<ul style="list-style-type: none"> <li>Most widely used data model</li> <li>Minimal complexity and provides a clear overview of the data</li> </ul>		<a href="#">MySQL, PostgreSQL - used to run queries and create reports e.g. consolidated customer statements</a>	
<b>Data Pipelines</b>	A set of data processing elements connected in series, where the output of one element is the input of the next one	<ul style="list-style-type: none"> <li>Enables the flow of data from an application to a data warehouse, from a data lake to an analytics warehouse, or into a system</li> </ul>		<a href="#">Data pipelines are used to perform predictive analysis e.g. a production department can use predictive analytics to know when the raw material is likely to run out, and it could also help forecast which supplier could cause delays</a>	
<b>ACORD Reference Architecture</b>	A series of seven interrelated industry models, or facets, which use different views to define the nature of the insurance industry	<ul style="list-style-type: none"> <li>Consists of: <ul style="list-style-type: none"> <li>Business glossary</li> <li>Informational models</li> <li>Data models</li> <li>Capability models</li> <li>Component models</li> <li>Process models</li> <li>Products frameworks</li> </ul> </li> <li>Implementation of ACORD standards improves data quality and improves data quality and flow, increase efficiency, and realize billion-dollar savings</li> </ul>		<a href="#">Official website</a>	

Term	Concept	Description /Analogy	Illustration	Real world Example	Relevance to Ki
<b>Batch-based data pipeline</b>	Data pipelines which are executed manually or recurringly	<ul style="list-style-type: none"> <li>In each run all data is extracted from the data source, operations are applied to the data and the processed data is published to a data sink</li> <li>Execution time ranges from a few minutes to a few hours depending on the quantity of data</li> </ul>		<a href="#">Payroll, billing, low frequency reports based on historical data</a>	
<b>Streaming data pipeline</b>	Data pipelines which are executed continuously all the time	<ul style="list-style-type: none"> <li>Consumes streams of messages, apply operations to the messages (such as transformations and filters) and publishes the processes messages to another stream</li> <li>Enables real-time processing</li> </ul>		<a href="#">Netflix uses streaming data pipelines to cater for its requirements of real time event-based processing and extensive support for customisation of windowing.</a>	
<b>Lambda architecture</b>	A data-processing architecture of three layers: batch processing, speed processing and a serving layer for responding to ad hoc queries	<ul style="list-style-type: none"> <li>Designed to handle massive quantities of data by taking advantage of both batch and stream processing methods</li> <li>Benefits include tolerance to human errors and hardware crashes, and scalability and quick response time</li> </ul>		<a href="#">Yahoo uses lambda architecture for running analytics on its advertising data warehouse</a>	
<b>Event-driven architecture</b>	A software architecture paradigm promoting the production, detection, consumption of, and reaction to events (an event is a change in state or an update e.g. an item being placed in a shopping cart)	<ul style="list-style-type: none"> <li>Uses events to trigger and communicate between decoupled services</li> <li>Has three key components: event producers, event routers and event consumers- producers publish an event to the router, which filters and pushes the event to consumers</li> <li>Common within modern applications built with microservices</li> </ul>		<a href="#">If a customer returns an item, this is logged by the stock management system (event producer), the event router pushes this to the finance system (event consumer) which updates to reflect this</a>	

Term	Concept	Description /Analogy	Illustration	Real world Example	Relevance to Ki
<b>Structured data</b>	Data that has been predefined and formatted to a set structure before being placed in data storage, which is often referred to as schema-on-write	<ul style="list-style-type: none"> <li>• Easily used by business users/ ML algorithms to draw business outcomes</li> <li>• Quick retrieval and less storage space required than unstructured data</li> </ul>		Credit card numbers, dates, address	
<b>Unstructured data</b>	Data stored in its native format and not processed until it is used	<ul style="list-style-type: none"> <li>• A variety of file formats that cannot be displayed in rows/columns as a relational database</li> <li>• Requires data science expertise/specialised tools (e.g. NLP) to extract value</li> </ul>		Emails, photos, social media posts, text, audio files	