

Data Strategy for AI (& Cloud)

Patipan Prasertsom

Slide developed in parts by Dr. Apivadee Piyatumrong, Assoc. Prof. Tiranee Achalakul, and Chayasin Saetia
with a lot of materials from AIGC-ETDA

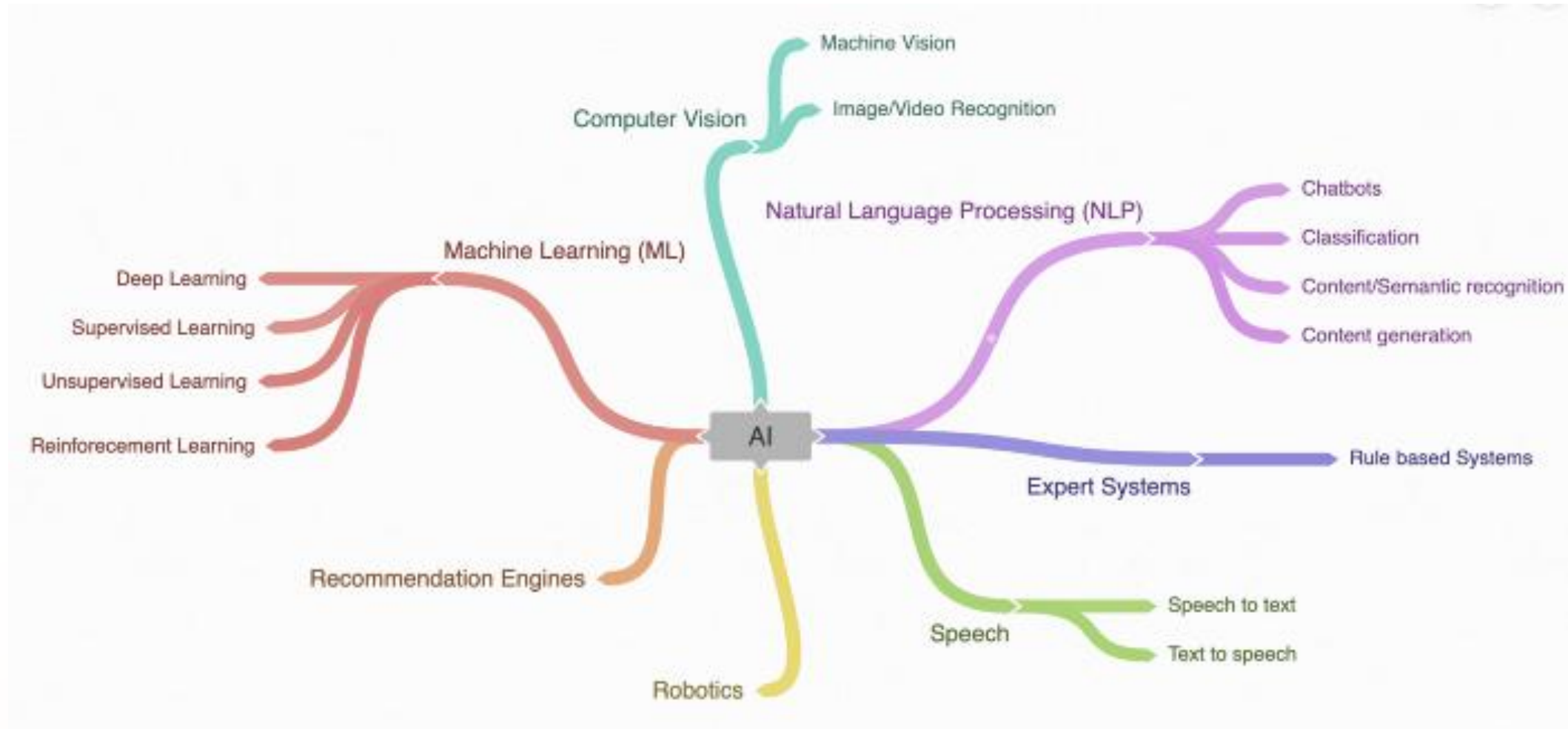
ARTIFICIAL INTELLIGENCE (A.I.)



AI is already in our life



AI is a broad field





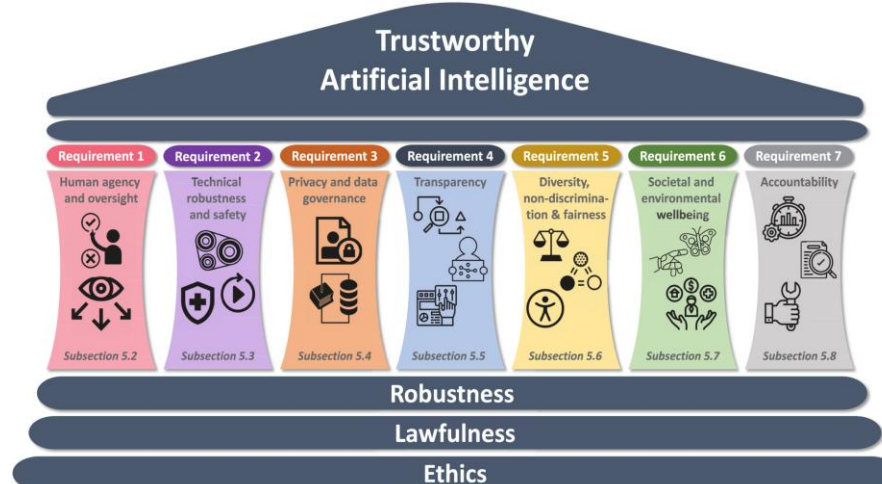
— GENERATIVE

AI

Generate new original content: text, images, videos and audio

- Generative AI learns to can create original utilizing patterns learned from training data
 - Content can be essays, solutions to problems, programming code, music score, realistic fakes, painting, and etc.
- Large Language Model (LLM) is a type of generative AI that excels in natural language tasks such as summarization, Q&A,etc.

GPT-3.5



2022

ChatGPT Launch

- Virality, "AI Boom"
- Mass adoption of conversational AI
- AI in education, coding, content

2023

Competition & Integration

- Major AI releases from competitors
- AI integration wave
- Point concerns and regulation

2024–2025

Multimodal & Personal AI

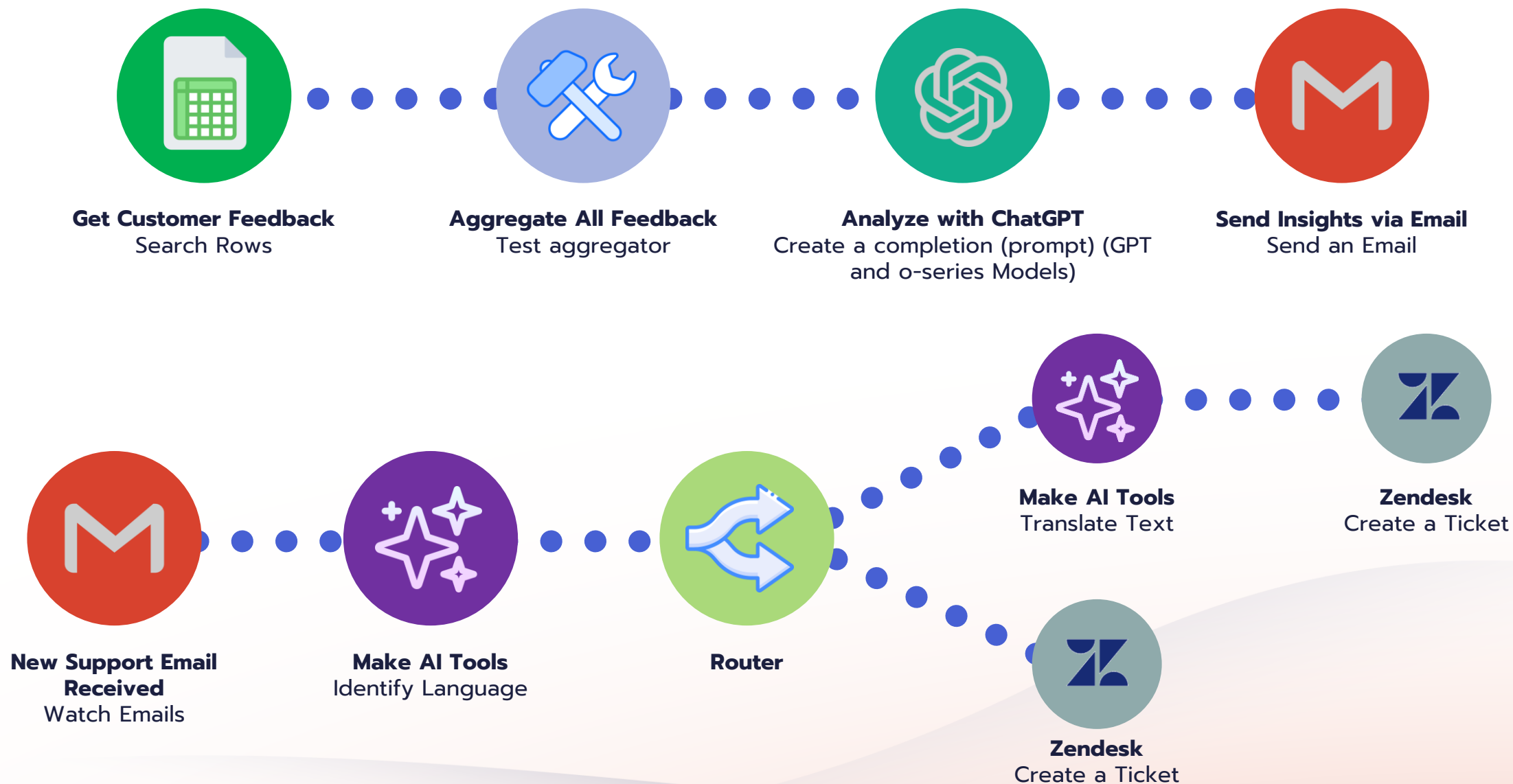
- Voice, image, video capabilities
- AI assistants and agents
- Reasoning and long-term memory

Personalized AI AI Automation Agentic AI

GPT-5

GPT-OSS-120b/20b (OpenW)
Claude Opus 4 / Sonnet 4
Gemini Deep Think
Llama 4 (not OpenW)
Gemma 3 (OpenW)
Magistral Small (reasoning)
QWEN 3

Examples of agentic workflows



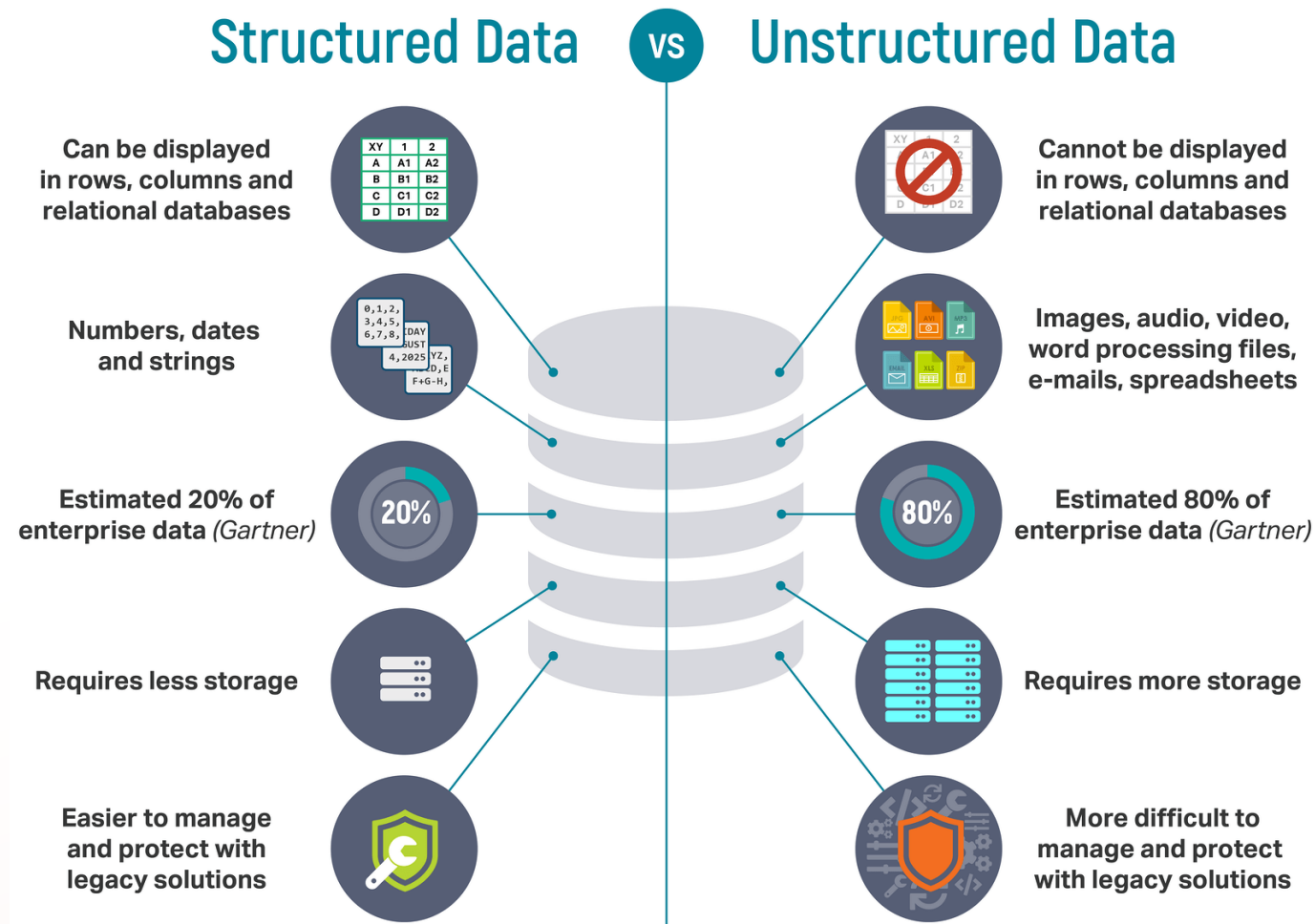
Data

(More on this later)

Mandatory Slide: Broad categories of collectible data

Suitable tasks:

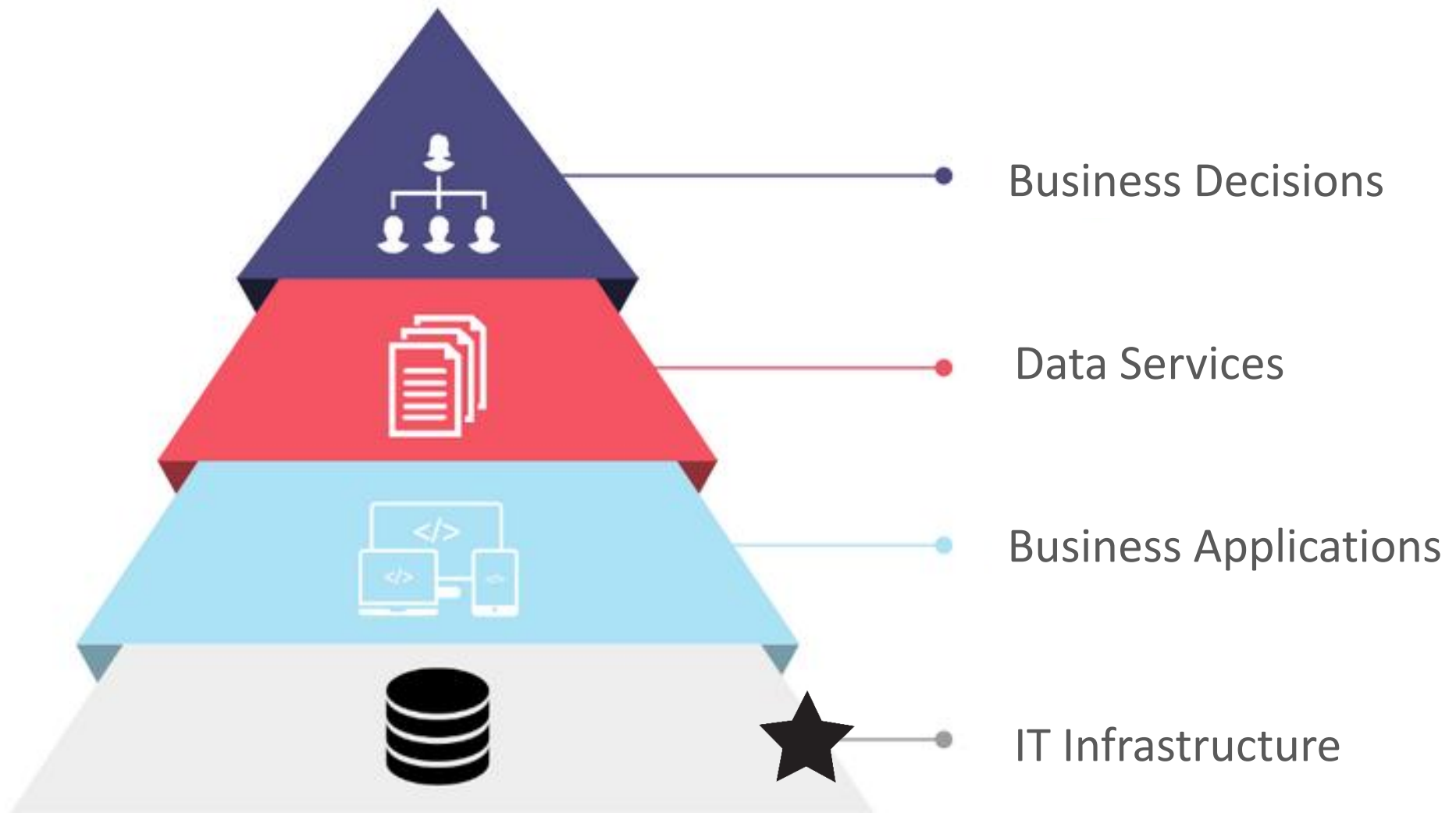
- Classical ML tasks
 - Prediction
 - Classification
 - Clustering
 - Etc.
- Data retrieval
- Analytics tasks



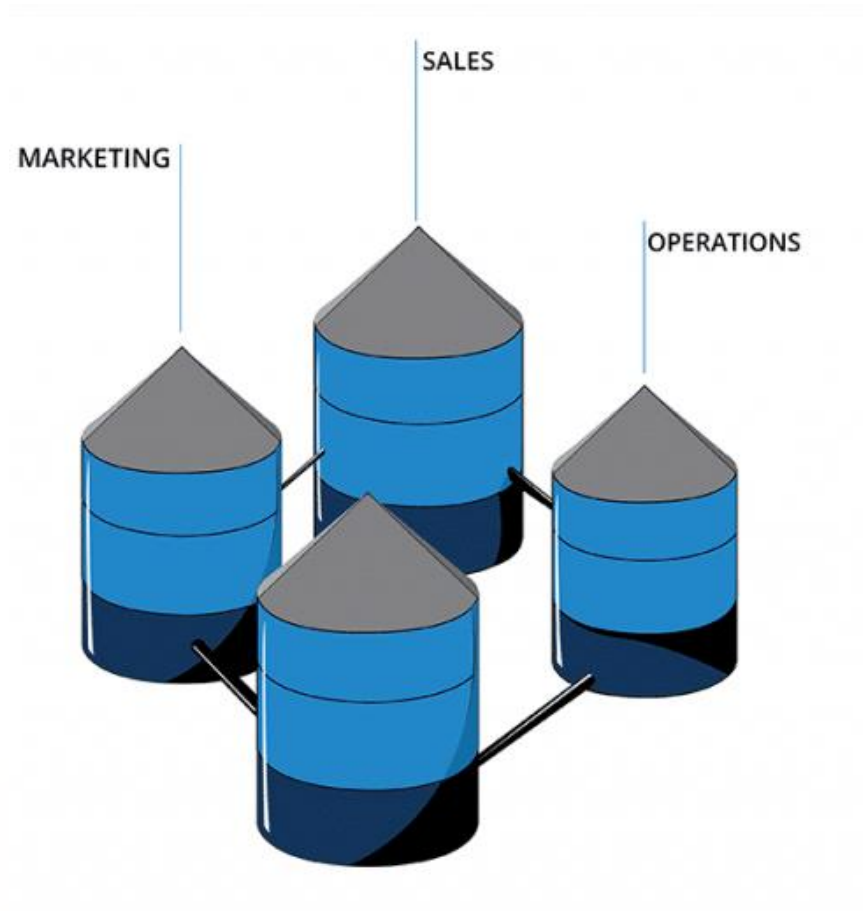
Suitable tasks:

- Natural Language Processing (NLP)
 - Text Analysis
 - Etc.
- Computer vision
- Generative AI and its applications

Strategy

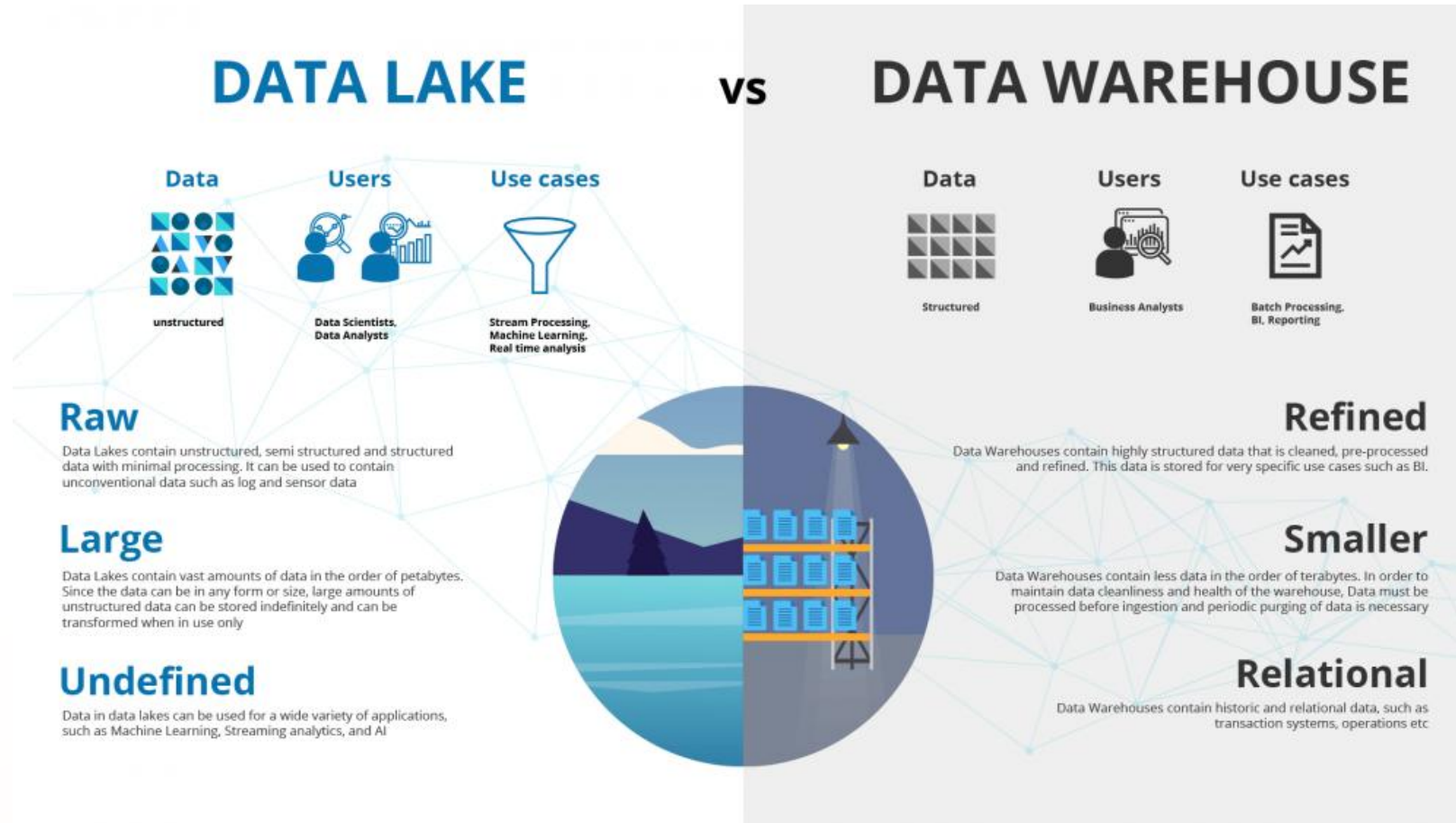


Common Issues: Data Silos



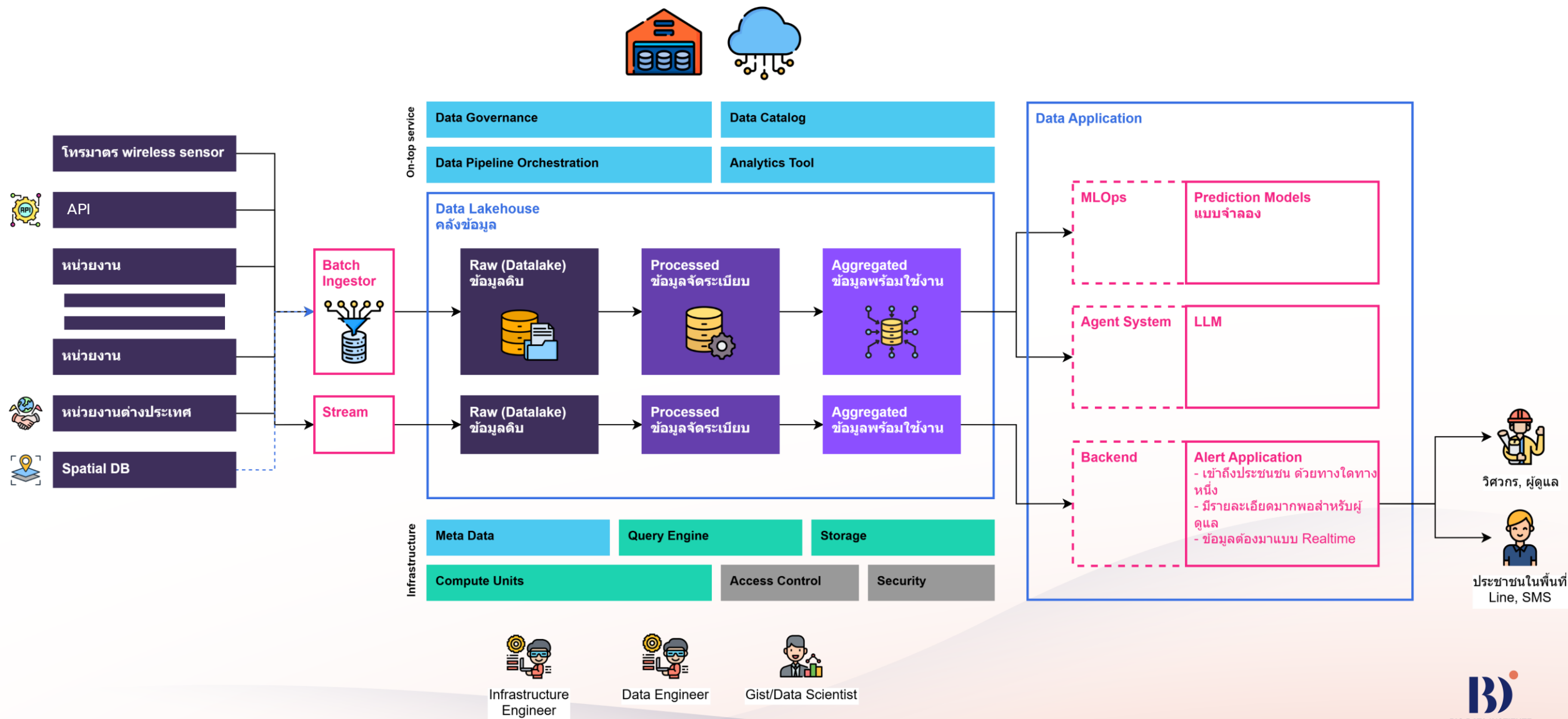
- Data Silos: separate data storages within an organization. Can result in issues such as:
 - Limited view of organization data
 - Data Inconsistency
 - Discourage collaborative work
- For AI/ML models to perform well, high quality data is needed (in preferably large amount)
 - There should be a centralized single source of truth for data.
 - This is also helpful for getting proper view of all available data

Common Solutions: Data Lake & Data Warehouse

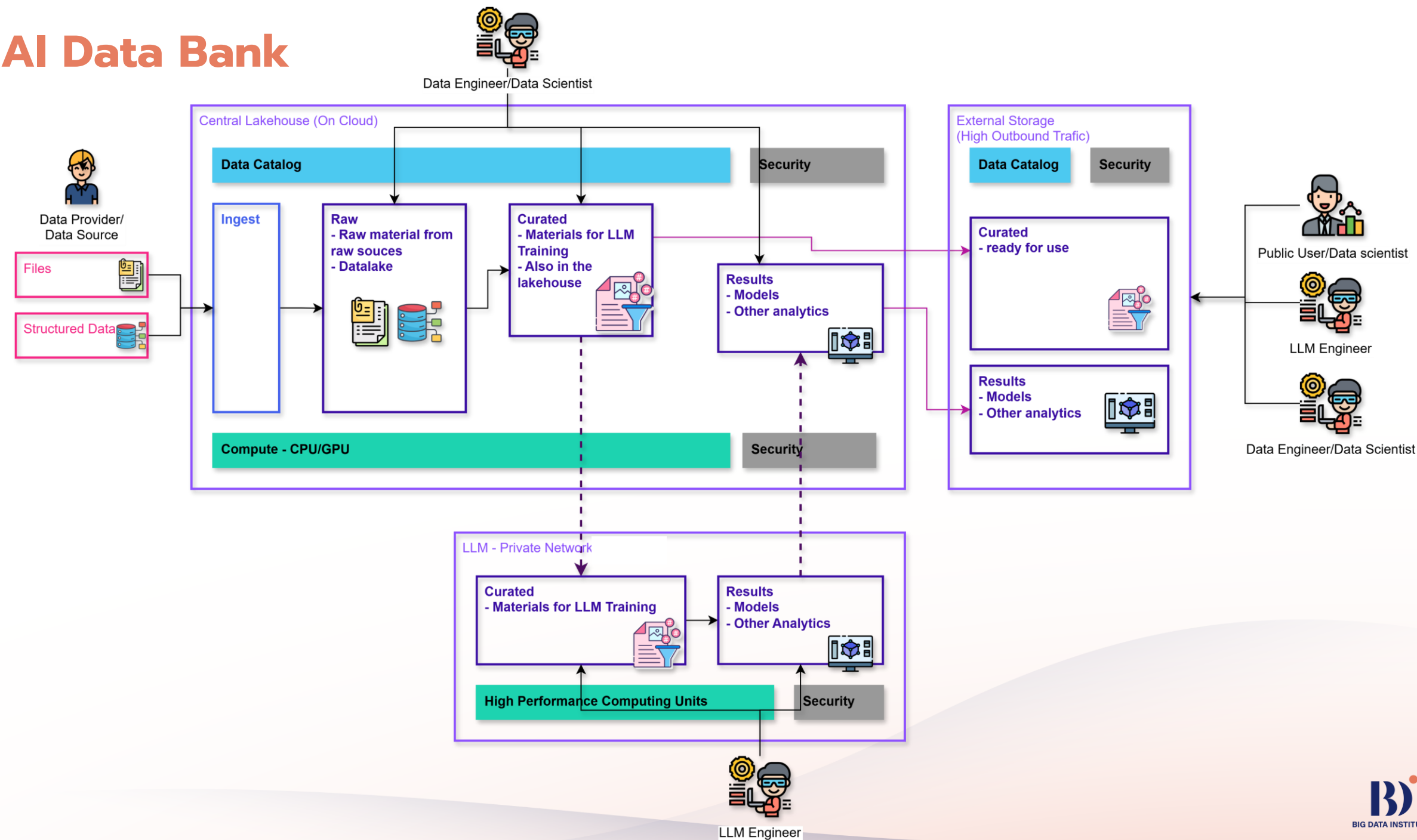


Nowadays, though, organizations have multiple different use cases, requiring both type of storages

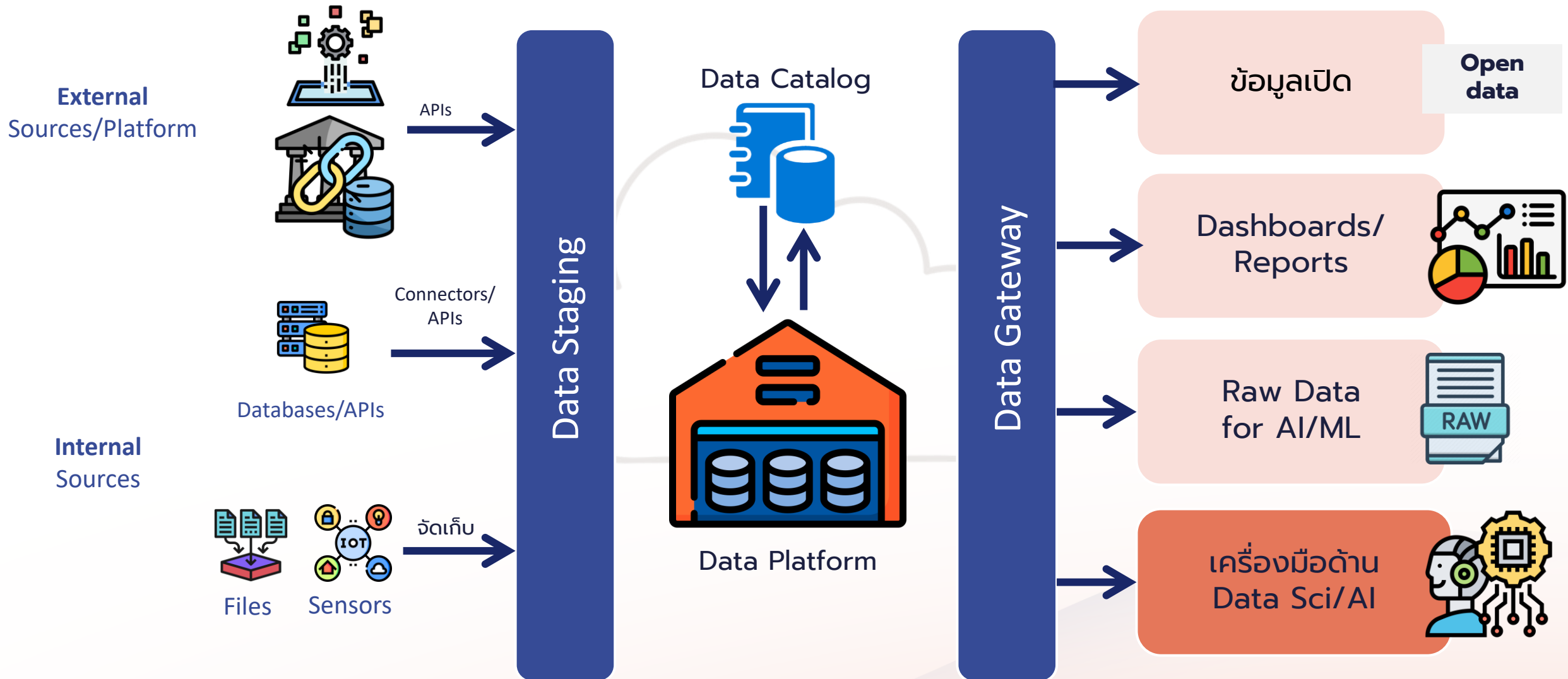
Data Lakehouse



AI Data Bank



Overall view of data platform



Choices for deployment



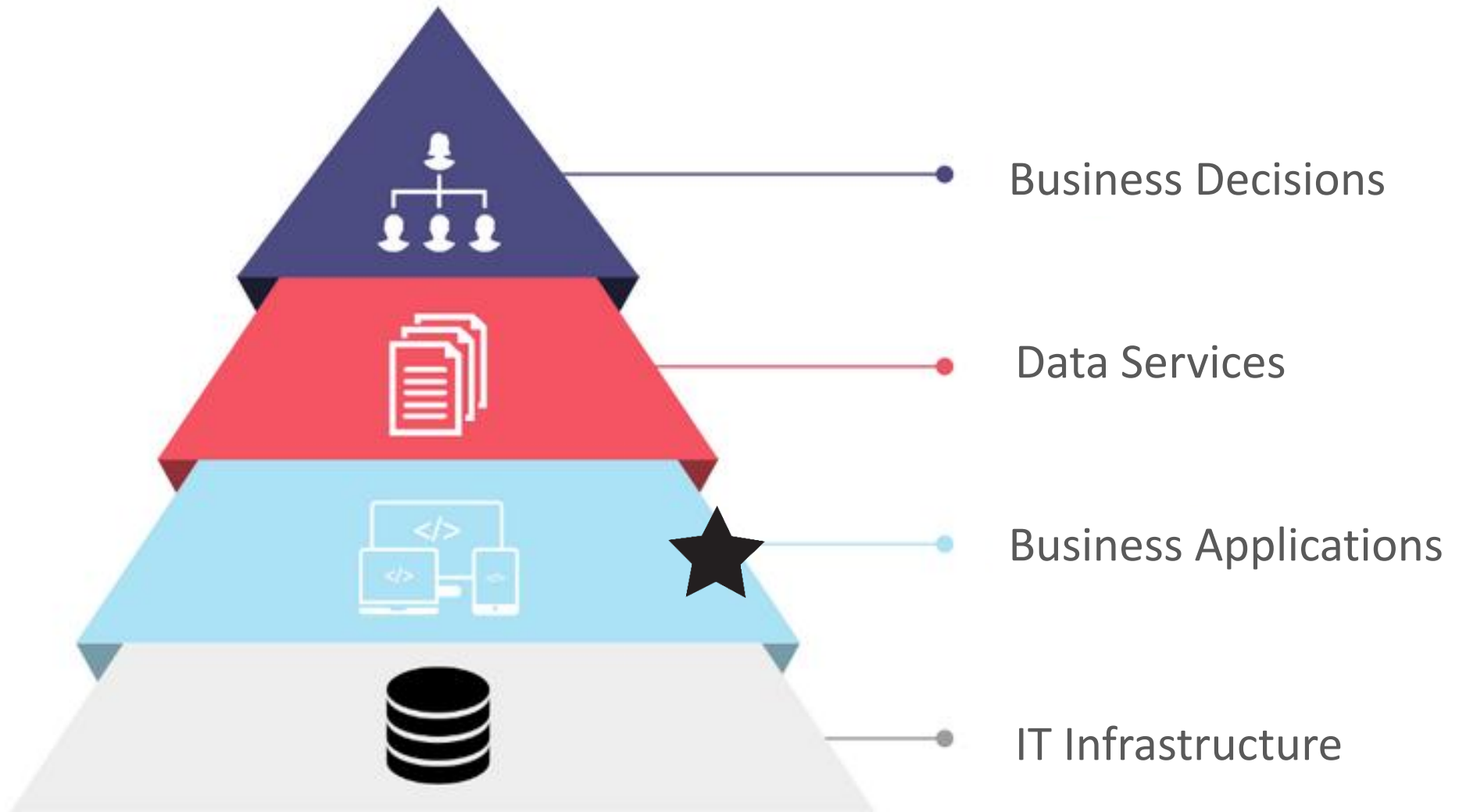
- **On-premise (local)**
 - Pros: Control over data/resources. Customizable.
 - Cons: High operating/maintenance costs
- **Cloud** Our topic
 - Pros: Low operation costs. Reliable. Scalable
 - Cons: Limited control. Require external connection
- **Hybrid**
 - Pros: Keep sensitive data locally. Cost saving
 - Cons: Complex to manage. High technical debts
- Need to decide
 - Which data goes on cloud?
 - How are we using the cloud?

Cloud service offerings

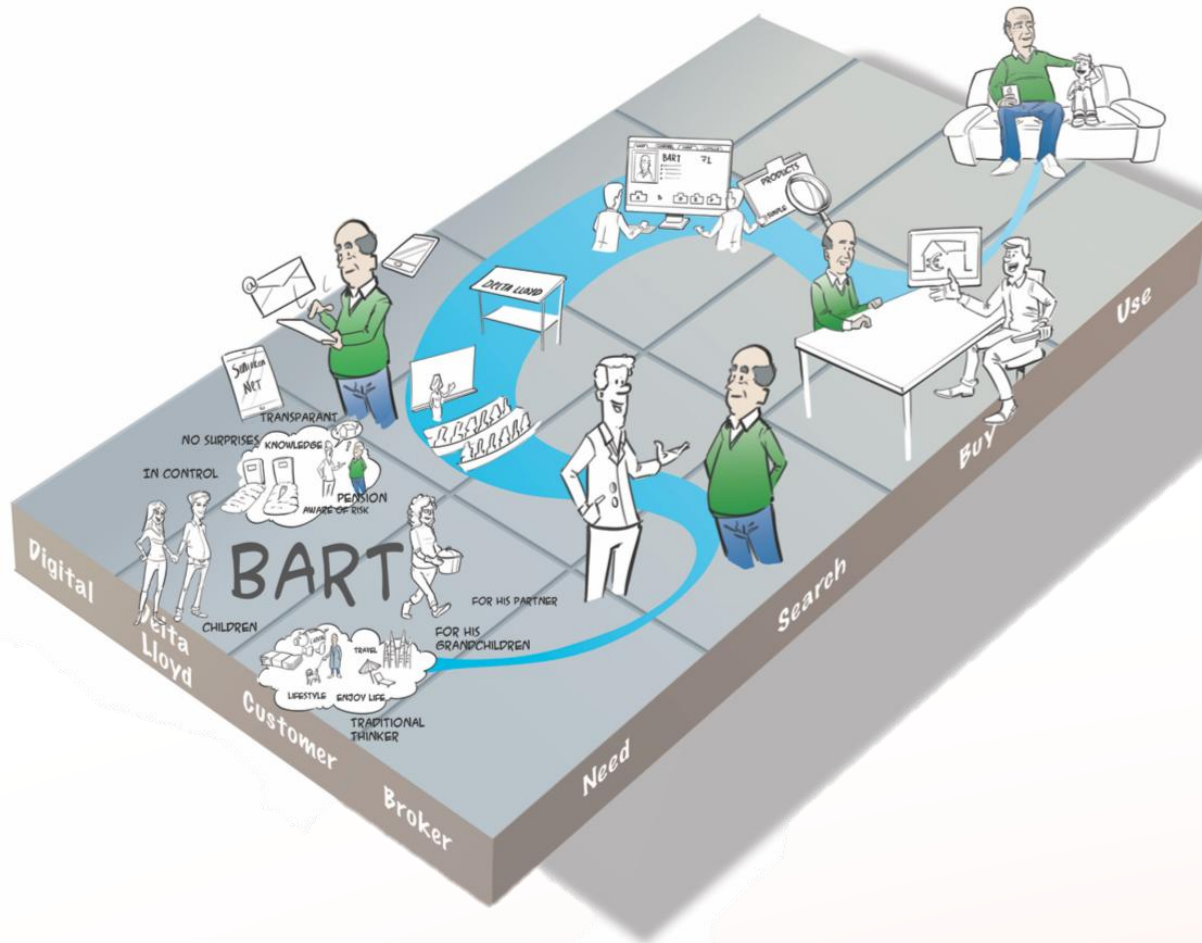
- **IaaS (Infrastructure as a Service)**
 - Provide “machines” (VM) as requested
 - Akins to lending private computers to users
- **PaaS (Platform as a Service)**
 - Provide “platform” to build applications/systems
 - Have pre-built components and tools available
- **SaaS (Software as a Service)**
 - Provide ready-to-use “software”
 - User can access complete program via internet

Not much to
talk about





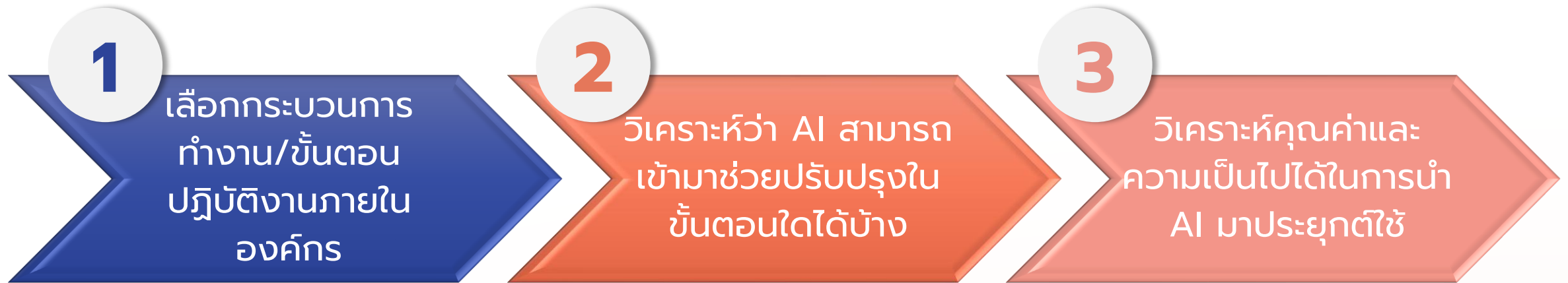
Business Journey

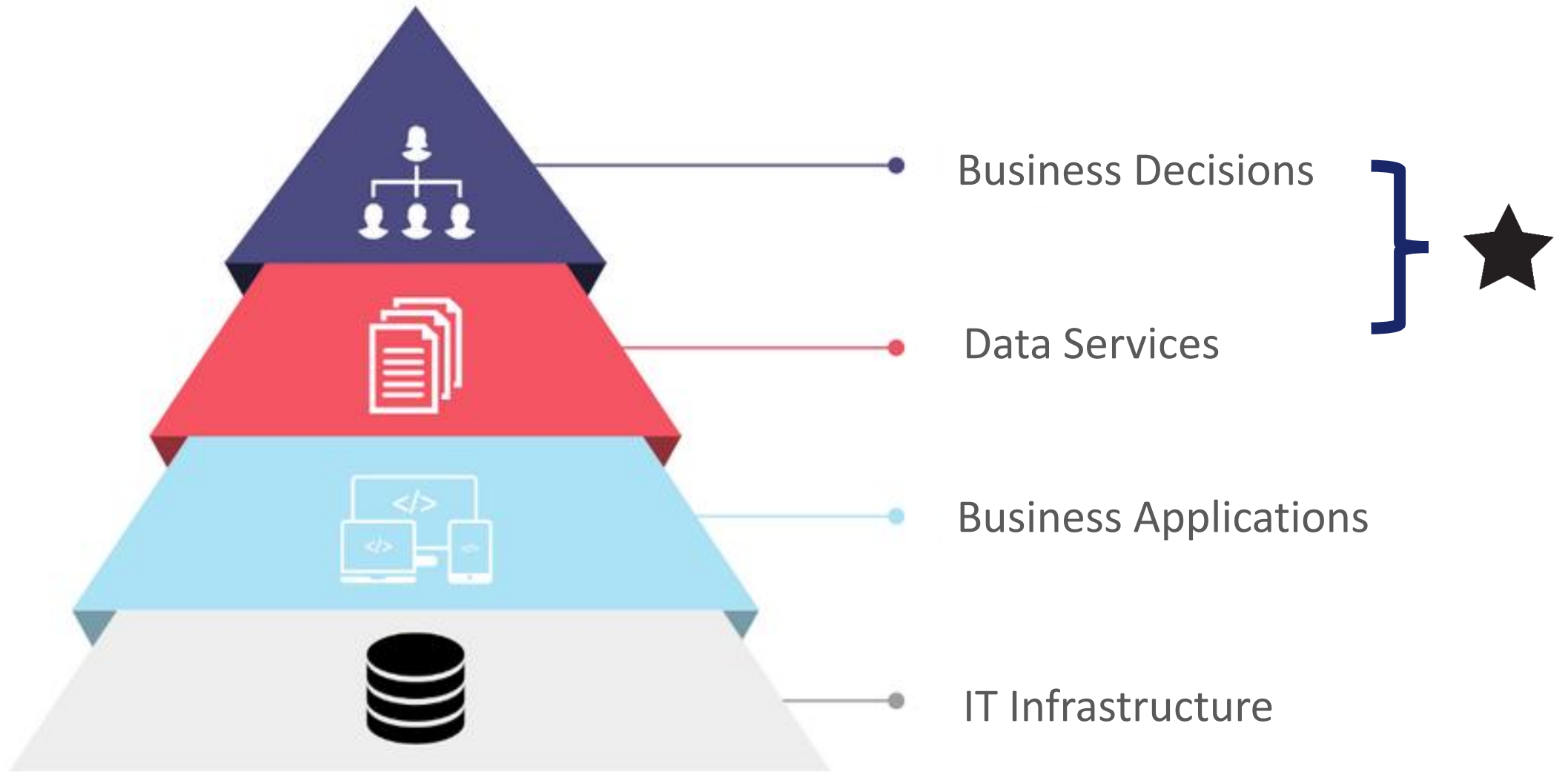


- It is almost never a good idea to blindly adopt technology
- Look for solutions to your problems, not the other way around
- Start from your pain points or things that should be improved
 - Look at journeys of your employees and customers

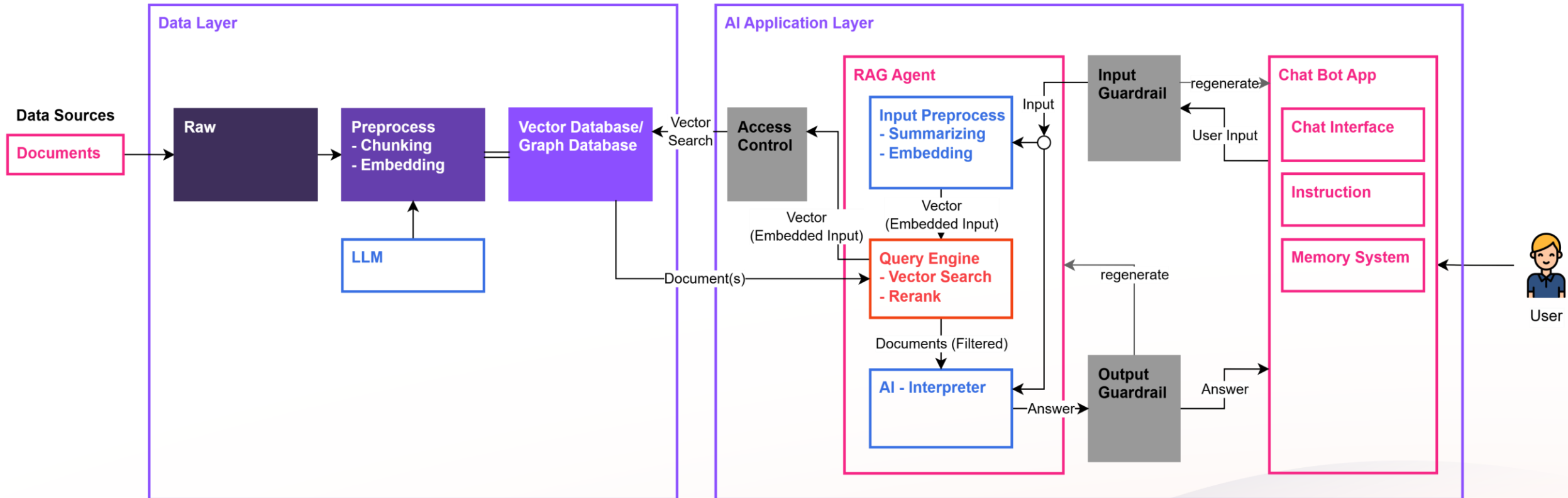
AI Use Cases Design

- ขั้นตอนการมองหาโอกาสในการนำ AI มาใช้งาน

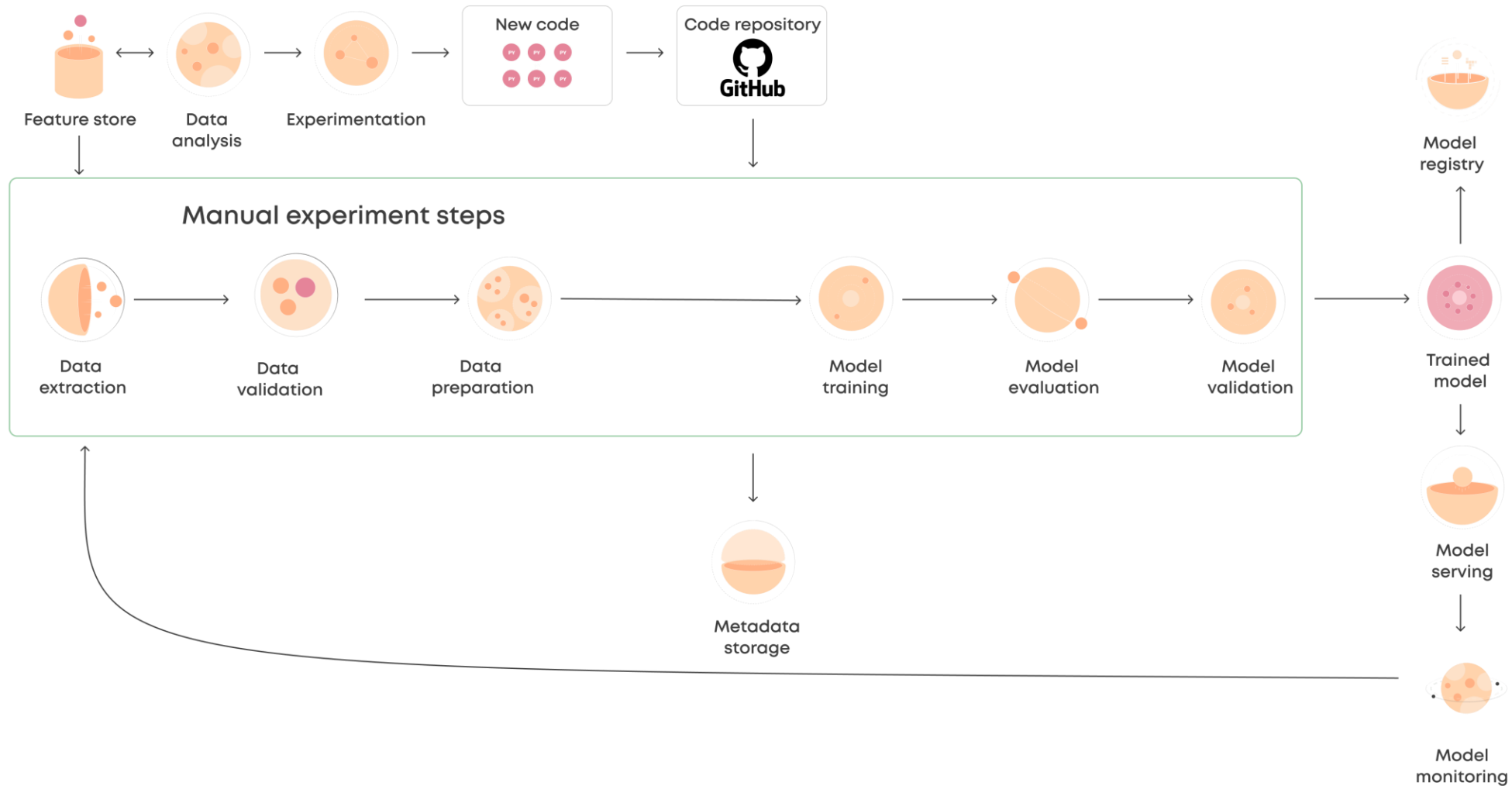




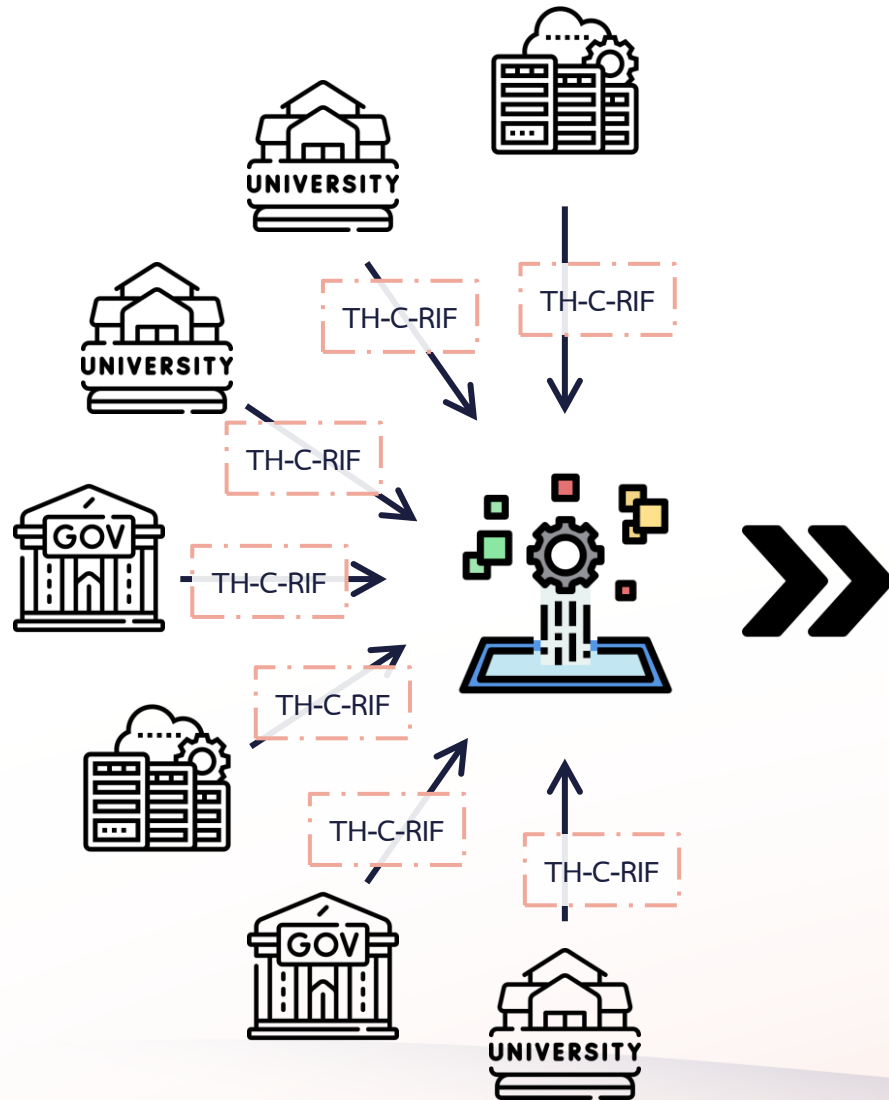
General RAG Pipeline (document-based)



General MLOps Pipeline



Planning for scale – Common Data Standard



USECASE 1

Research Performance and Reporting ประกอบการกำหนดนโยบายและยุทธศาสตร์ ววน.

USECASE 2

Project Approval, Project Monitoring and Evaluation สำหรับการบริหารจัดการการให้ทุน

USECASE 3

One-Stop-Shop Research Information Service เพื่อให้บริการข้อมูลววน.แบบเชื่อมโยงได้

USECASE อื่น ๆ

การวิเคราะห์ข้อมูลโดย Interactive Dashboard และแบบจำลองทางคณิตศาสตร์

Data Governance

Data Governance : ธรรมาภิบาลข้อมูล

Accuracy Completeness Timeliness Security Privacy Connectedness Worthiness

Organization & Stewardship

- Function of
- DG Committee
 - Steward team
 - Data Controller/Processor/User

Data management Policy

- Data Lifecycle
- Data Security & Privacy
- Data Quality Assurance
- Data Exchange

Audit

- Data Risk Management
- Law & Regulation Compliance
- Data Quality Audit

Building Knowledge & Awareness

- Program Coverage
- Measurement
- Ongoing

Data Management Policy : นโยบายการบริหารจัดการข้อมูล

Data Lifecycle

Create

Store

Use

Archive

Destroy

Data Catalog

- What data ?
- Who's the owner ?
- Search Tags
- Data sources

Data Security & Privacy

- Confidentiality
- Availability
- Integrity

Data Quality

- Accuracy
- Validity
- Timeliness
- Completeness
- Uniqueness
- Consistency

Foundation

Data Exchange

Data Governance

- Essentially, we want to be able to answer:
 - What data ?
 - Who owns it ?
 - How often is it updated?
 - Where's it from ?
- Along with being able to properly manage them



Security Control



Quality Control



Access Control

Data Classification

		หลักการจำแนกข้อมูล				
หมวดหมู่ข้อมูล (ภาครัฐ)	หมวดหมู่ข้อมูล (ภาครัฐ)	เปิดเผยได้ (OPEN)	ใช้ภายใน (PRIVATE)	ลับ (Confidential)	ลับมาก (Secret)	ลับที่สุด (Top Secret)
	ตัวอย่างข้อมูล	รายงานศึกษาทางวิชาการ	ข้อมูลการปฏิบัติงาน	ข้อมูลสัญญาที่รออนุมัติ	ข้อมูลพนักงานรายบุคคล	ข้อมูลกำลังรบ
	ข้อมูลเปิด	✓ ข้อมูลเปิด ภาครัฐ				
	ข้อมูลใช้ภายใน		✓ ใช้หลักเกณฑ์ที่มีผลกระทบต่อองค์กรในการประเมินและจำแนก			
	ข้อมูลส่วนบุคคล		✓ อิงตาม พ.ร.บ.คุ้มครองข้อมูลส่วนบุคคล พ.ศ. 2562			
	ข้อมูลข่าวสารลับ			✓ พ.ร.บ. ข้อมูลข่าวสารของทางราชการ พ.ศ. 2540 ระเบียบว่าด้วยการรักษาความลับของทางราชการ พ.ศ. 2544		
	ข้อมูลความมั่นคง			✓ นโยบายและแผนระดับชาติ ว่าด้วยความมั่นคงแห่งชาติ (พ.ศ 2562 – 2565)		

What about using/storing data on cloud?

Data Classification (2)

37

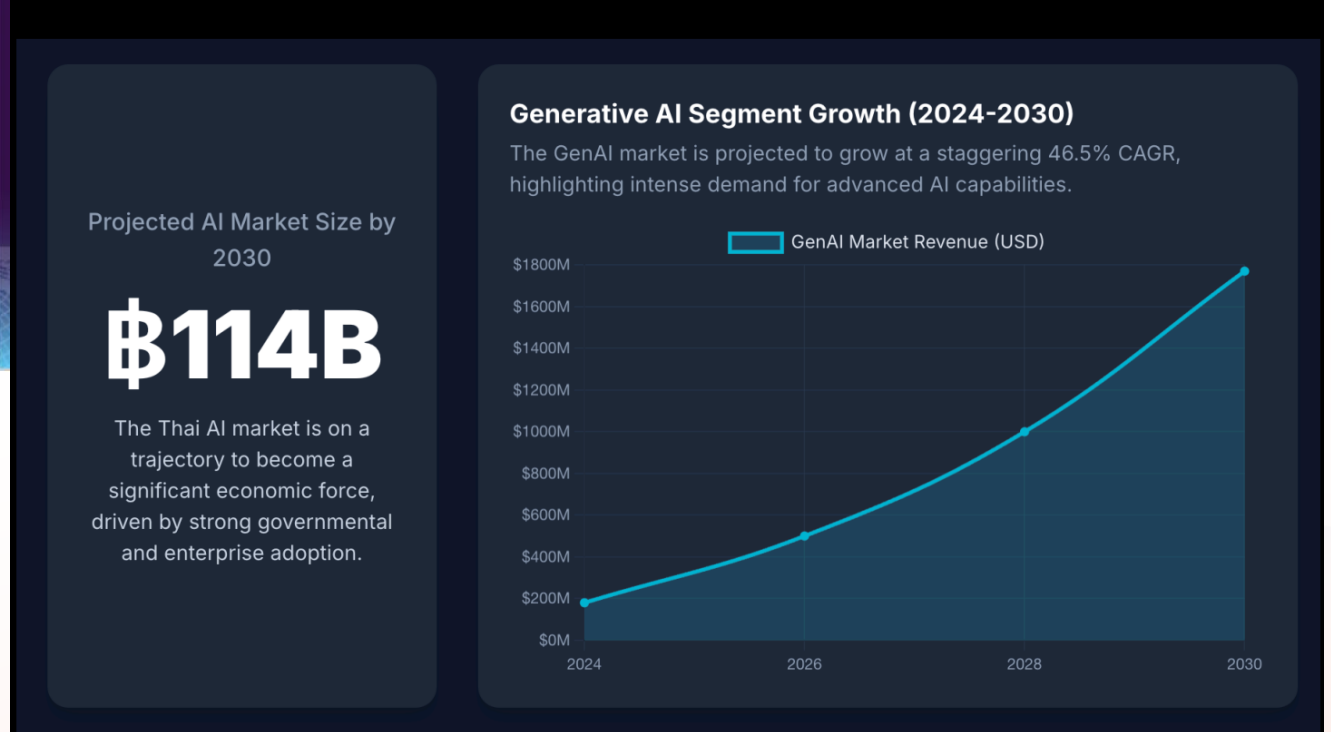
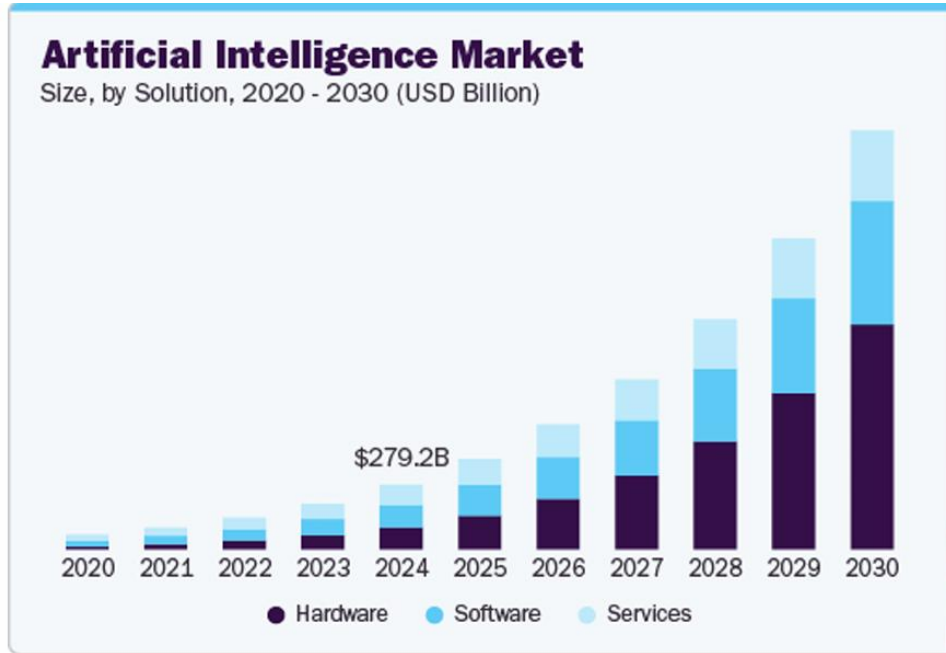


	ระดับความอ่อนไหวของข้อมูล (Sensitivity Level)	การเข้าถึง (Accessibility) (SAMPLE)	การจำแนกข้อมูล (Classification Label)	ตัวอย่างข้อมูล (Sample Type of Data)
1	Level 0 - None	ทุกคนสามารถเข้าถึงได้	เปิดเผยได้ (Public Used)	<ul style="list-style-type: none"> ข้อมูล PR ข้อมูล Website ข้อมูลบริษัท
2	Level 1 - Low	พนักงานทั่วไป	ใช้ภายใน (Internal Used Only)	<ul style="list-style-type: none"> ข้อมูลบริการ ข้อมูลการตลาด ข้อมูลจัดซื้อ
3	Level 2 - Medium	ระดับผู้บริหาร / IT / คณะทำงาน กรรมการบริหารข้อมูล (ต้องมีการขอสักสิทธิ์ในการเข้าถึงข้อมูล)	ลับ (Confidential)	<ul style="list-style-type: none"> ข้อมูลการเงิน ข้อมูลรายได้การขาย ข้อมูล NDA / สัญญา
4	Level 3 - High	ระดับผู้บริหาร / IT / คณะทำงาน กรรมการบริหารข้อมูล (ต้องมีการขอสักสิทธิ์ในการเข้าถึงข้อมูล)	ลับมาก (Highly Confidential)	<ul style="list-style-type: none"> ข้อมูลลูกค้า ข้อมูลภายใน IT System ข้อมูลพนักงานรายบุคคล
5	Level 4 - Extreme	เจ้าของข้อมูล (Data Owner) (ต้องมีการขอสักสิทธิ์ในการเข้าถึงข้อมูล)	ลับที่สุด (Restricted)	<ul style="list-style-type: none"> ข้อมูลส่วนบุคคล Credit Card ข้อมูลส่วนบุคคลที่เป็นข้อมูลอ่อนไหว

13

Data Governance For AI

Why do we particularly care about AI?



AI Governance



Competitiveness &
Sustainability



Transparency &
Accountability



Compliance with Laws,
Ethics, Standard

Human-centric
decision
making



Security & Privacy



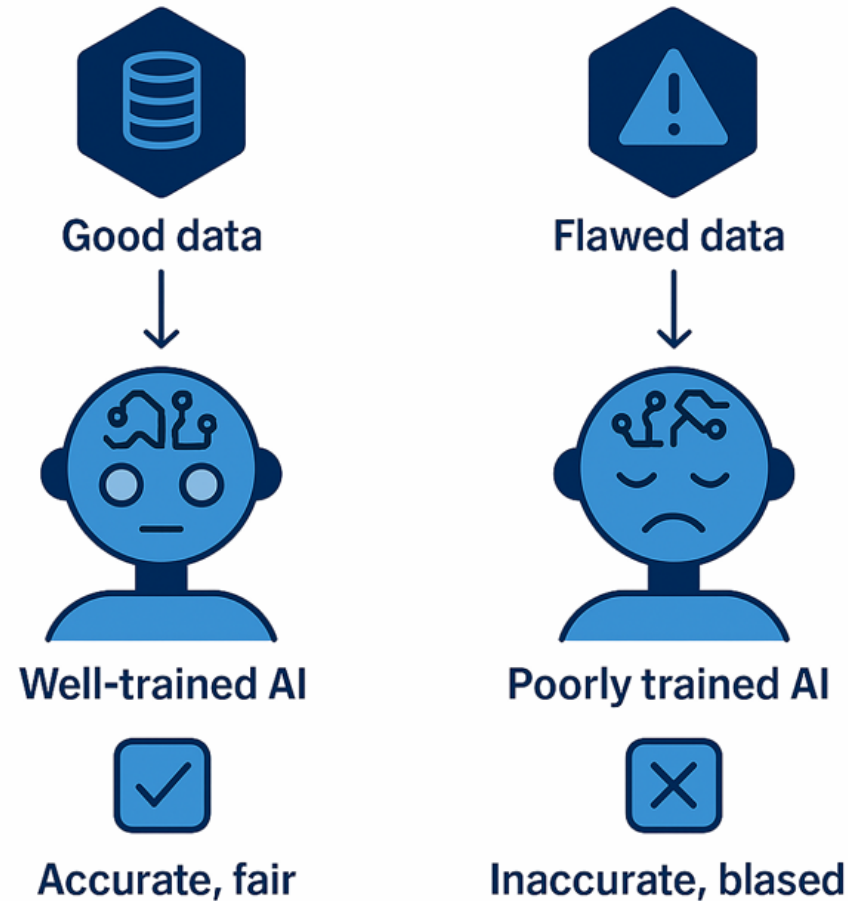
Equality & Fairness



Reliability

What role do data play in AI Governance?

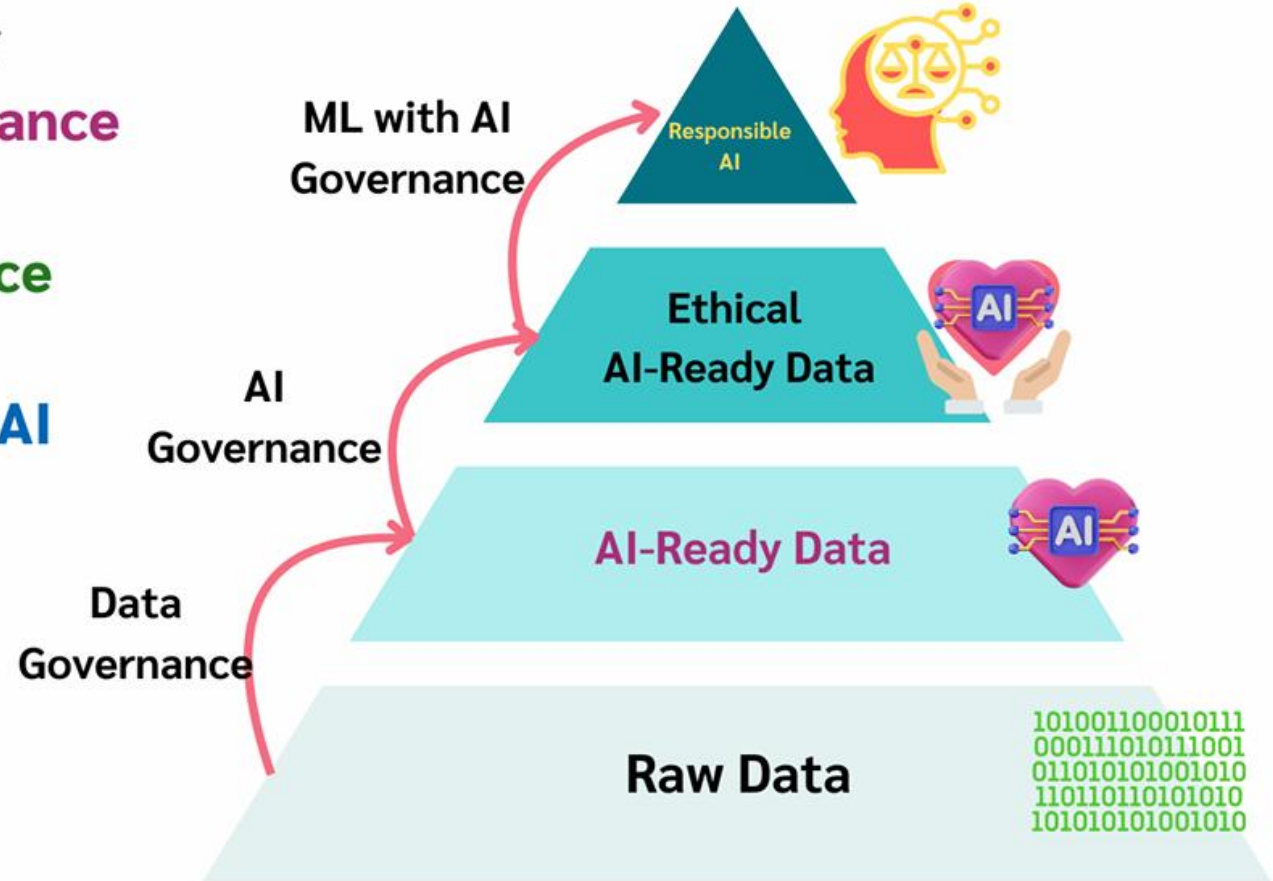
เมื่อระบบ AI ได้รับ ข้อมูลที่มีคุณภาพ (Good Data) เช่น ข้อมูลที่ถูกต้อง ครบถ้วน และไม่ลำเอียง ส่งผลให้ AI ที่ถูกฝึกมีความสามารถในการวิเคราะห์ได้อย่างถูกต้อง (Accurate) และ เป็นธรรม (Fair)



เมื่อระบบ AI ได้รับ ข้อมูลที่บกพร่อง (Flawed Data) เช่น ข้อมูลที่ลำเอียง ขาดสมดุล หรือมีอคติ ส่งผลให้โมเดล AI มีแนวโน้ม ตัดสินใจผิดพลาด (Inaccurate) และ ไม่เป็นธรรม (Biased) ซึ่งอาจก่อให้เกิดผลกระทบทั้งเชิงนโยบาย เศรษฐกิจ และสิทธิมนุษยชน

Data Governance for AI

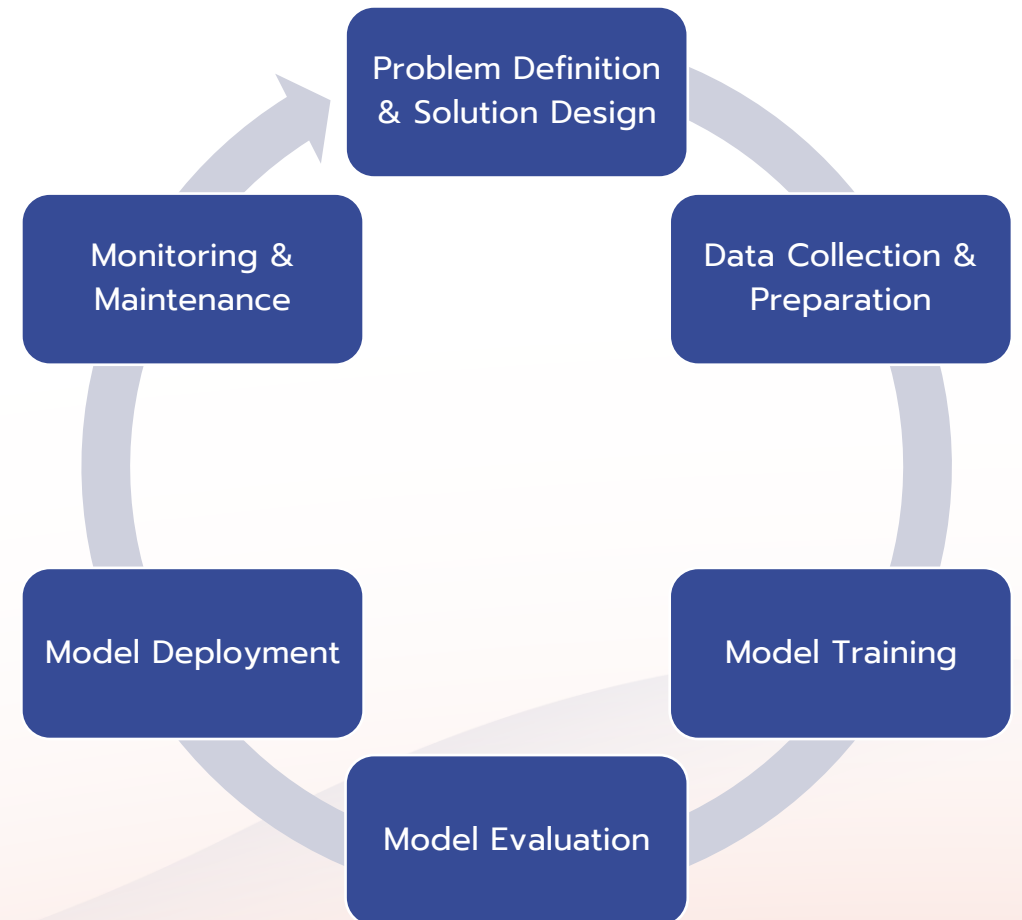
การบูรณาการ
Data Governance
และ
AI Governance
เพื่อมุ่งสู่
Responsible AI



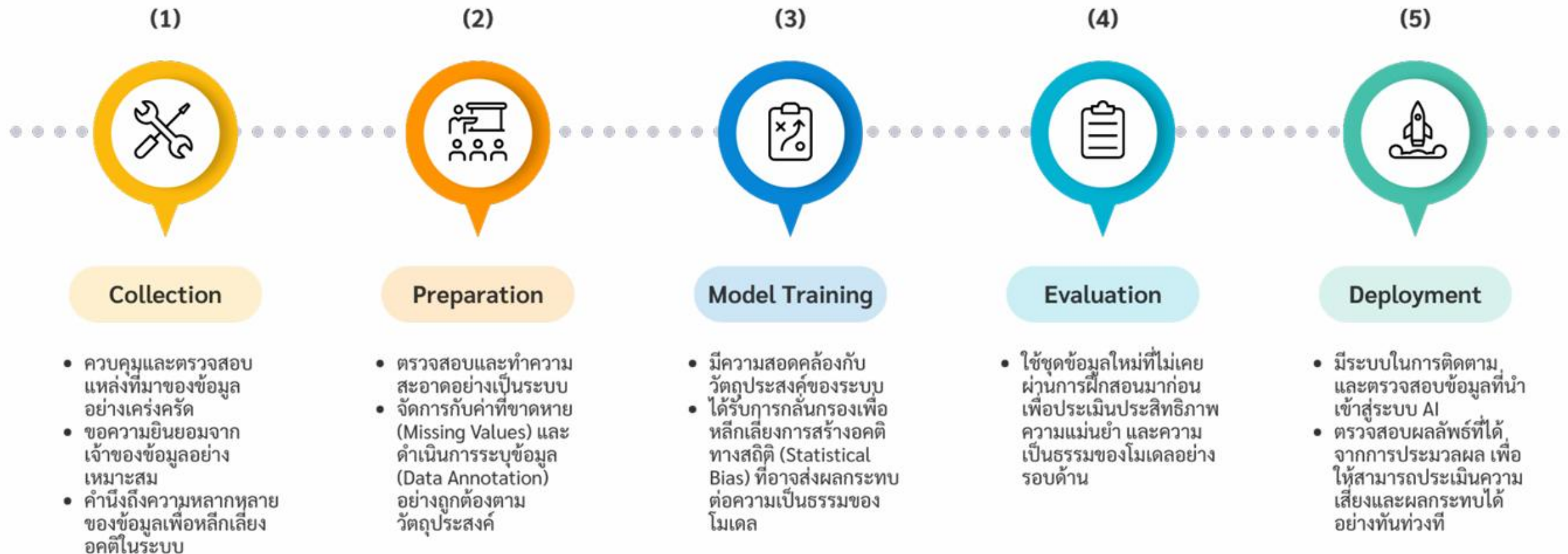
AI-Ready Data & AI Lifecycle



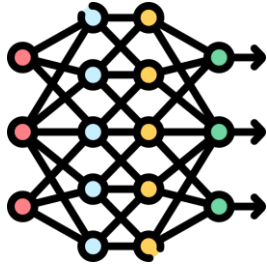
AI Lifecycle



Bringing them together



Data Preparation Guideline



Original Model (from scratch)

- Use diverse datasets that covers all required scenario
- Perform **quality control**
- Contains **metadata** and **data lineage**
- Accurate Labeling of data



Open Source/ Pretrained Model

- Investigate **inherent bias**
- Use data that **reflect actual usage** for fine-tuning or testing
 - Transparent & Traceable
 - Do not drastically differ from original training data
 - Appropriate privacy handling
- Perform **fairness evaluation**



Off-the-shelf (service/software)

- Consider **Input appropriateness**
 - Sensitive data?
 - Personal data?
- **Validate outputs** before actual deployment
- Perform **risk assessment**

Data Quality Criteria

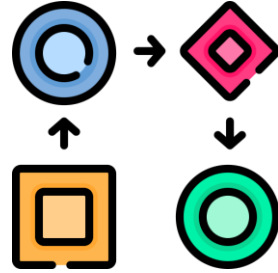
เกณฑ์	ความหมาย
Accuracy	ข้อมูลถูกต้อง เป็นไปตามความจริง
Consistency	ข้อมูลมีความสอดคล้อง ไม่ขัดแย้งกันเอง
Availability	สามารถเข้าถึง/ใช้งานข้อมูลได้ทั้งในปัจจุบันและในอนาคต
Completeness	ข้อมูลมีความสมบูรณ์ ครบถ้วน
Conformance	ข้อมูลเป็นไปตามมาตรฐานที่เป็นที่ยอมรับ
Credibility	ข้อมูลได้มาจากแหล่งข้อมูลที่เชื่อถือได้
Processability	ข้อมูลสามารถถูกใช้งานและประมวลผลได้โดยระบบสารสนเทศ
Relevance	ข้อมูลมีจำนวนมากพอที่จะใช้งานได้อย่างมีประสิทธิภาพ
Timeliness	ข้อมูลมีความทันสมัย ใหม่พอที่จะใช้ประโยชน์ได้อย่างมีประสิทธิภาพ
Context	ข้อมูลเหมาะสมที่จะใช้งานในบริบทที่ต้องการ

4 Key Components of Data Governance for AI



Roles & Responsibilities

- Data Owner
- Data Steward
- Data Custodian
- AI Model Owner



Metadata & Data Lineage

- **Metadata:** data that explain the context of data
 - Format, owners, update freq., etc.
- **Data Lineage:** the process of tracking data over time
 - Origin, end point, transformation



Policy & Standards

- Purpose Limitation
- Consent Management
- Data Management Policy/Lifecycle



Risk Management

- Data Bias
- Privacy Risk
- Legal Risk
- Mitigation:
 - Bias Assessment
 - De-identification
 - Impact Assessment

Example of Roles & Responsibility in AI Life Cycle

RACI CHART

AI Life Cycle
(Roles with Responsibilities)

Roles	Collect	Prepare	Train	Monitor
Data Scientist	C	C	R	R
Data Engineer	R	R	A	C
Data Analyst	C	I	C	I
Project Manager	A	A	C	A

- R: Responsible
- A: Accountable
- C: Consulted
- I: Informed

AI-Ready Data Governance Framework

Layer 1: Data Layer – Foundation for AI-Ready Data



Data Quality

Accurate, complete, up-to-date, and consistent data



Systematic Metadata

Documentation of what, where, and when data is used



Data Lineage

Traceable processes prior to AI model use



Transparent Collection

Defined purpose and user consent for data collection

Layer 2: Governance Layer – Systematic Oversight



Roles & Responsibilities

Clear designation of responsible parties

Policy & Standards

Principles for data usage and ethics impact



Access & Consent Management

Restricted data access and proper consents

Risk & Impact Assessment

Analysis of technical, ethical, and social risks

Layer 3: Ethics Layer – Alignment with Global Standards



Fairness

No bias against particular groups



Transparency

Disclosure of AI decision-making



Explainability

Reasonable justifications for outcomes



Privacy

Protection of personal data

Thank You