# Efficient training of machine learning algorithms
## — Optimisation of results at reduced costs

Heiko Joerg Schick
Chief Architect & Industry Expert | Advanced Computing, Artificial Intelligence & Semiconductor

Presenting the work of many people at Huawei

**2024-03-19 v2 | Generative AI DACH 2024, Berlin, Germany**

HUAWEI TECHNOLOGIES DÜSSELDORF GmbH

HUAWEI

# Agenda

**Huawei Pangu-Weather**

– Data and settings

– Computational costs

**Implementation of end-to-end lifecycle in AI projects**

**Fine-tuning**

– GPT assistant training pipeline | LLMs model sizes over time | Full-parameter fine-tuning

– How is Low-Rank Adaption (LoRA) different?

– Number of trainable parameters | Percent of total parameters

– QLoRA

**Retrieval-augmented generation (RAG) system**

– Overcoming challenges of LLMs

– General Purpose AI (GPAI) classification and key requirements for providers

– Fine Tuning vs. Retrieval Augmented Generation

– Basic chatbot architecture | Example

– Retrieval-augmented generation (RAG) architecture | Example | Many questions ?!?

**Closing remarks**

**Heiko Joerg Schick**
Chief Architect & Industry Expert
*Advanced Computing, Artificial Intelligence & Semiconductor*
Munich Research Center

**HUAWEI TECHNOLOGIES**
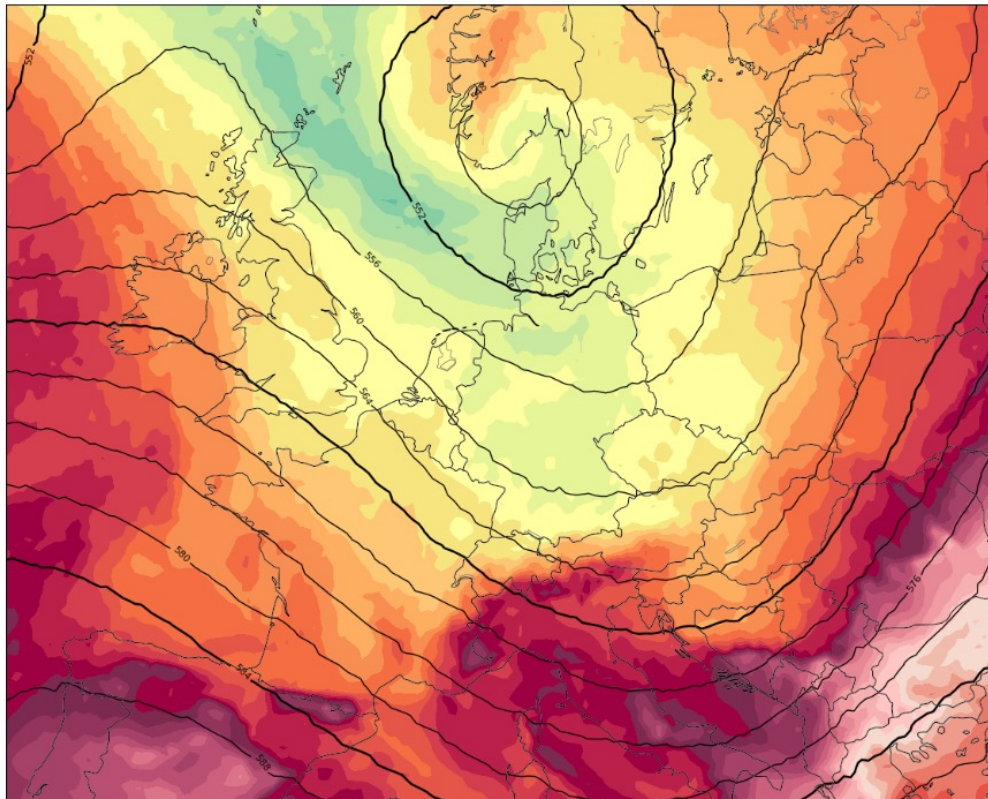**Duesseldorf GmbH**
Riesstrasse 25, D0
80992 Munich

Mobile +49-151-54682218
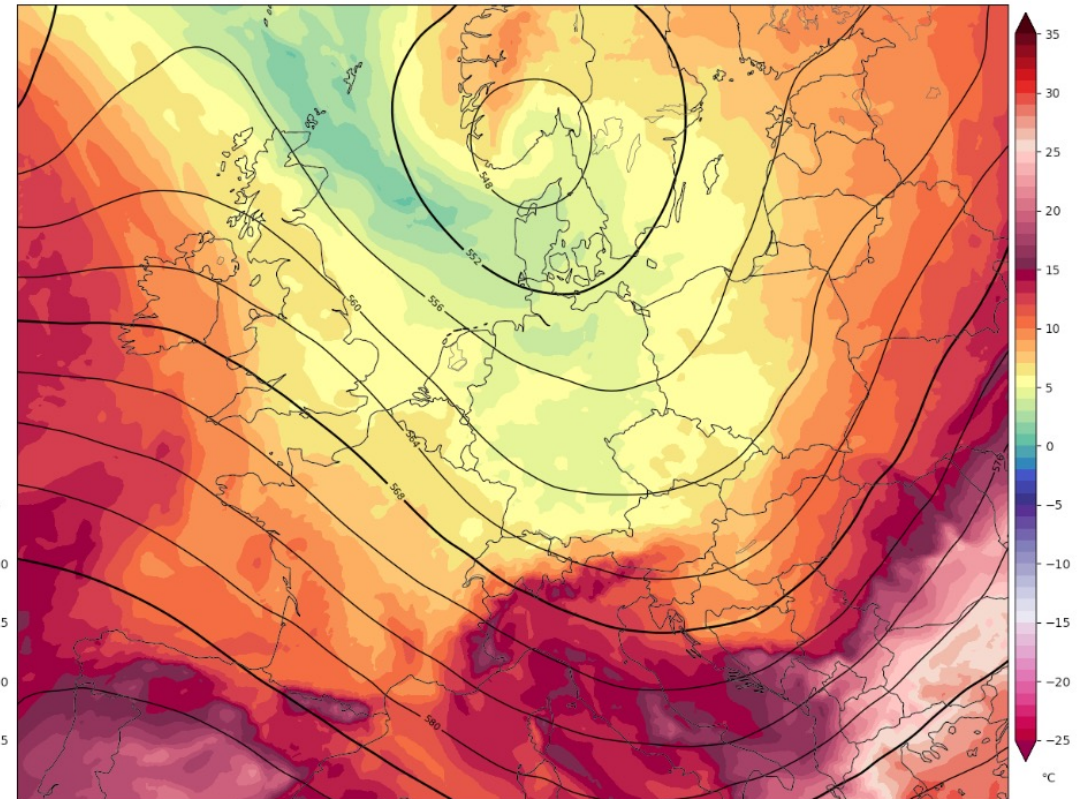E-mail heiko.schick@huawei.com

LinkedIn https://www.linkedin.com/in/heikojoergschick/

# Huawei Pangu-Weather – ICON Comparison Maps (Bi et al., 2023)



Pangu-Weather (initialized from ICON-Analysis)

ICON-Global-Deterministic

Shaded: Temperature at 850hPa, Lines: Geopotential Height at 500hPa in gpdm
Run: Wed, 26 Jul 2023, 12UTC, Valid Date: Wed, 26 Jul 2023, 12UTC (+0h)

With the current setup, a single 7-day forecast with Huawei Pangu-Weather consumes 14 Wh of energy. For a 7-day forecast with the ICON model, the energy consumption amounts to approximately 30000 Wh. This simple calculation of course does not include the energy consumption required to generate the training data and to train the model.

# Huawei Pangu-Weather (Bi et al., 2023)
— Data and settings

- The dataset includes the **5th generation of ECMWF reanalysis (ERA5) data**, which is publicly available.

- It comprises hourly reanalysis data from the year 1940 onwards.

- For our study, we used data from **1979 to 2017 for training** purposes, **2019 data for validation**, and **2018, 2020, and 2021 data for testing** to ensure a fair comparison with WeatherBench.

- The dataset contains a variety of surface and upper-air variables across **37 pressure levels**.

- Specifically, we selected **four surface variables** (2m temperature, u- and v-components of 10m wind speed, mean sea-level pressure) and **five upper-air variables** (geopotential, specific humidity, temperature, u- and v-components of wind speed) at **13 selected pressure levels** (ranging from 50hPa to 1000hPa).

- Although the **full dataset exceeds 2000 TB** in size, our analysis used approximately **60 TB** of data.

# Huawei Pangu-Weather (Bi et al., 2023)
## — Computational costs

- The training phase involves each forecast model having approximately **64 million parameters**.

- Each model is trained for **100 epochs** over **16 days** using 192 NVIDIA Tesla V100 GPUs, indicating that the models have not yet converged.

- During inference, each forecast takes about **1.4 seconds** on a single V100 GPU.

- Inference can also be carried out on a CPU, albeit with a longer processing time.

- Executing a 7-day global forecast involves running the 24-hour model seven times, totalling **less than 10 seconds**.

- Faster inference facilitates easier ensemble forecasting.

# Implementation of end-to-end lifecycle in AI projects (Alake, 2020), (Sato et al., 2019)

- Problem statement
- Ideal problem solution
- Understanding and insight into the problem
- Technical requirements

- Data structure and source
- Solution form
- Model architecture
- Algorithm research
- Hardware requirements

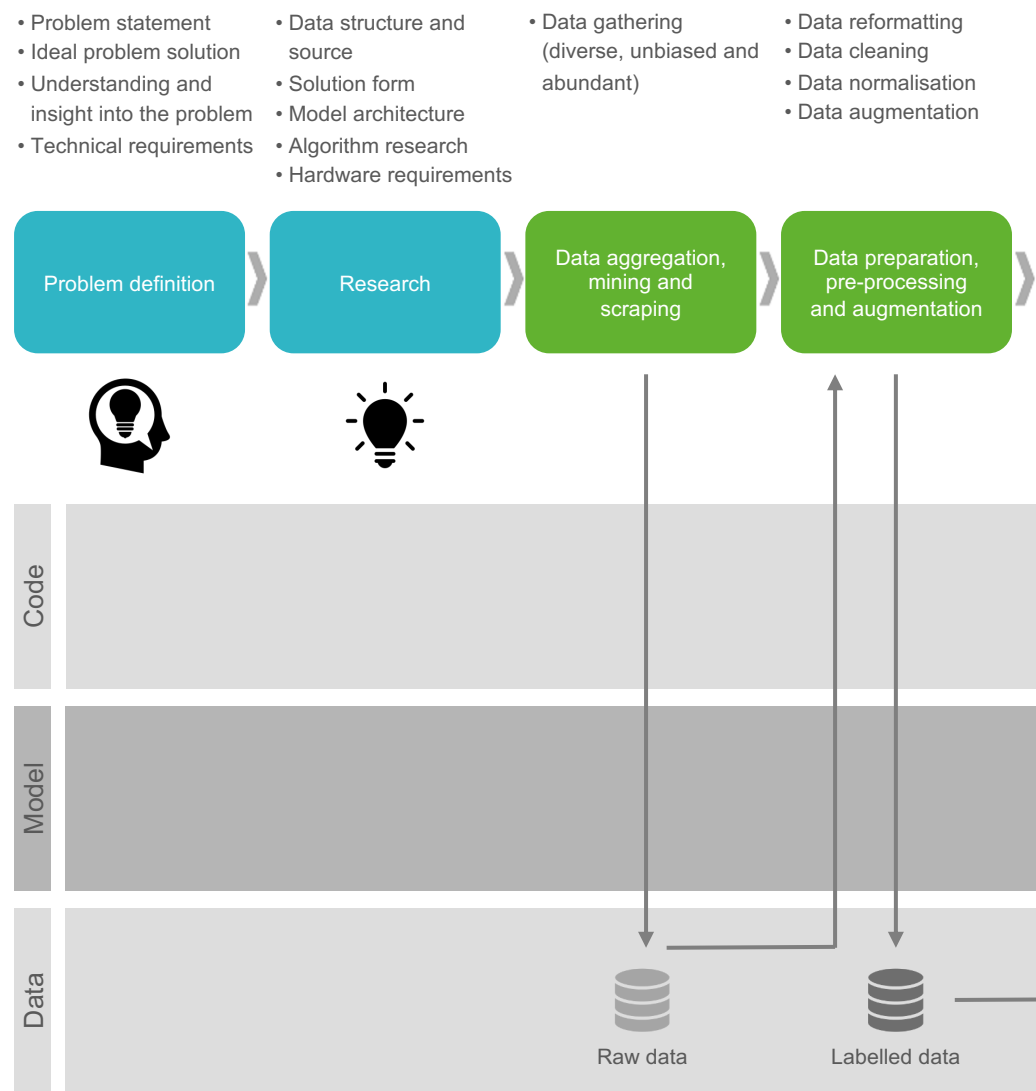Problem definition 〉 Research 〉

Code

Model

Data

# Implementation of end-to-end lifecycle in AI projects (Alake, 2020), (Sato et al., 2019)

- Problem statement
- Ideal problem solution
- Understanding and insight into the problem
- Technical requirements

- Data structure and source
- Solution form
- Model architecture
- Algorithm research
- Hardware requirements

- Data gathering (diverse, unbiased and abundant)

- Data reformatting
- Data cleaning
- Data normalisation
- Data augmentation

| Problem definition | Research | Data aggregation, mining and scraping | Data preparation, pre-processing and augmentation |

Code

Model

Data

Raw data     Labelled data

# Implementation of end-to-end lifecycle in AI projects (Alake, 2020), (Sato et al., 2019)

- Problem statement
- Ideal problem solution
- Understanding and insight into the problem
- Technical requirements

- Data structure and source
- Solution form
- Model architecture
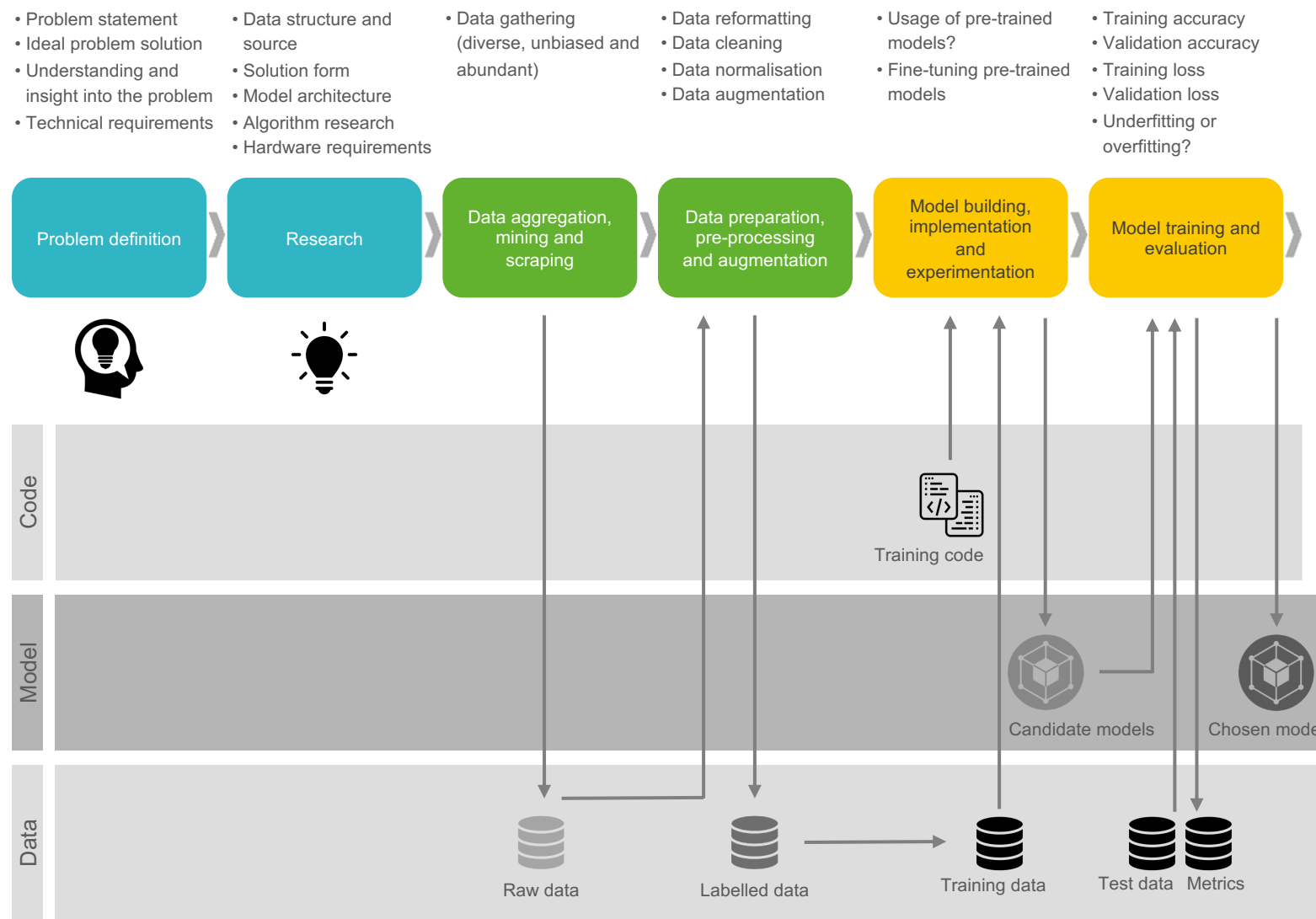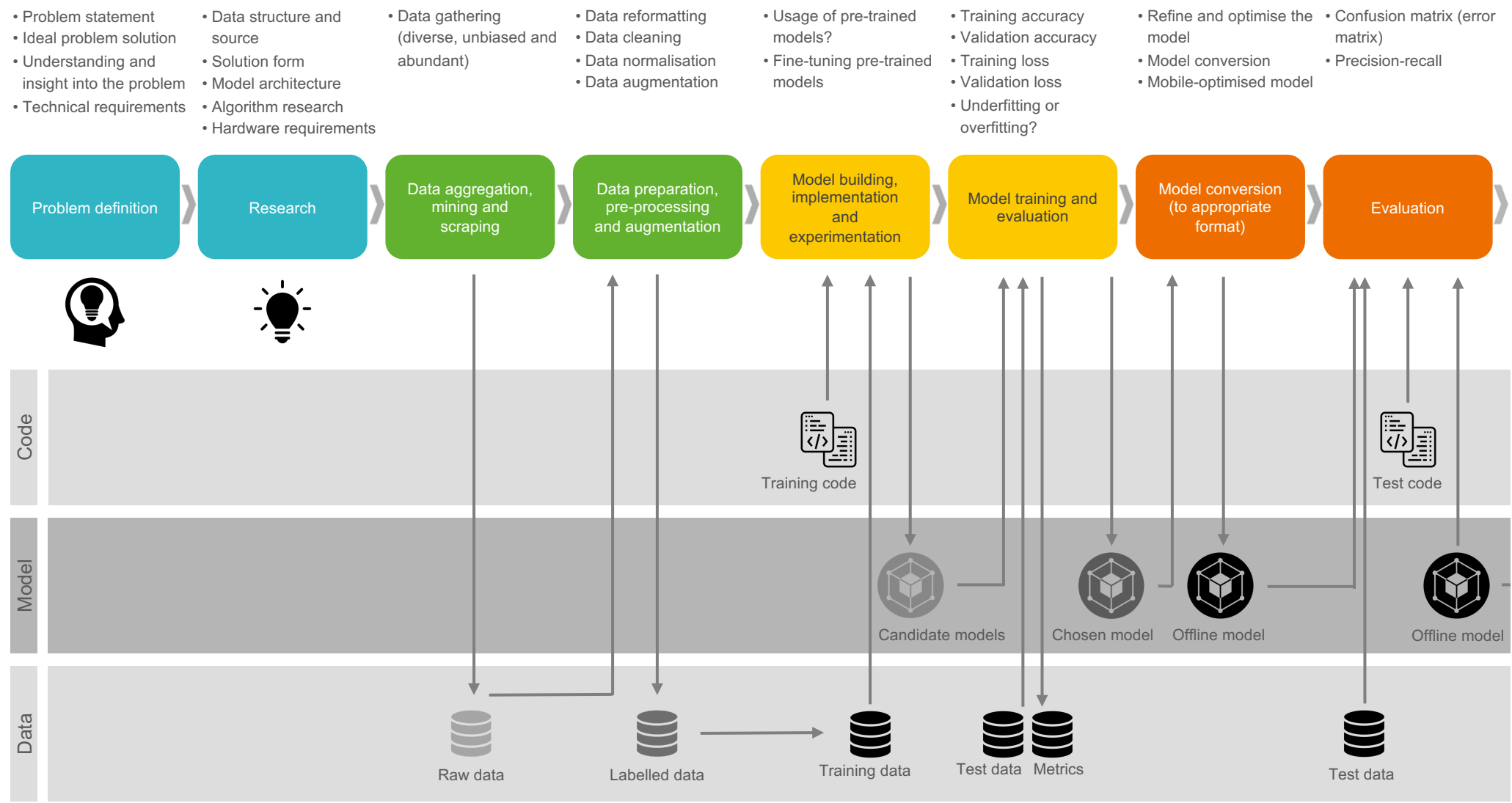- Algorithm research
- Hardware requirements

- Data gathering (diverse, unbiased and abundant)

- Data reformatting
- Data cleaning
- Data normalisation
- Data augmentation

- Usage of pre-trained models?
- Fine-tuning pre-trained models

- Training accuracy
- Validation accuracy
- Training loss
- Validation loss
- Underfitting or overfitting?

| Problem definition | Research | Data aggregation, mining and scraping | Data preparation, pre-processing and augmentation | Model building, implementation and experimentation | Model training and evaluation |

**Code**

Training code

**Model**

Candidate models    Chosen model

**Data**

Raw data    Labelled data    Training data    Test data    Metrics

# Implementation of end-to-end lifecycle in AI projects (Alake, 2020), (Sato et al., 2019)

- Problem statement
- Ideal problem solution
- Understanding and insight into the problem
- Technical requirements

- Data structure and source
- Solution form
- Model architecture
- Algorithm research
- Hardware requirements

- Data gathering (diverse, unbiased and abundant)

- Data reformatting
- Data cleaning
- Data normalisation
- Data augmentation

- Usage of pre-trained models?
- Fine-tuning pre-trained models

- Training accuracy
- Validation accuracy
- Training loss
- Validation loss
- Underfitting or overfitting?

- Refine and optimise the model
- Model conversion
- Mobile-optimised model

- Confusion matrix (error matrix)
- Precision-recall

| Problem definition | Research | Data aggregation, mining and scraping | Data preparation, pre-processing and augmentation | Model building, implementation and experimentation | Model training and evaluation | Model conversion (to appropriate format) | Evaluation |

**Code**

Training code

Test code

**Model**

Candidate models

Chosen model

Offline model

Offline model

**Data**

Raw data

Labelled data

Training data

Test data   Metrics

Test data

# Implementation of end-to-end lifecycle in AI projects (Alake, 2020), (Sato et al., 2019)

- Problem statement
- Ideal problem solution
- Understanding and insight into the problem
- Technical requirements

- Data structure and source
- Solution form
- Model architecture
- Algorithm research
- Hardware requirements

- Data gathering (diverse, unbiased and abundant)

- Data reformatting
- Data cleaning
- Data normalisation
- Data augmentation

- Usage of pre-trained models?
- Fine-tuning pre-trained models

- Training accuracy
- Validation accuracy
- Training loss
- Validation loss
- Underfitting or overfitting?

- Refine and optimise the model
- Model conversion
- Mobile-optimised model

- Confusion matrix (error matrix)
- Precision-recall

- UI interface to access model functionalities
- Continuous integration pipeline that enables model redeployment

- Model performance monitoring system

| Problem definition | Research | Data aggregation, mining and scraping | Data preparation, pre-processing and augmentation | Model building, implementation and experimentation | Model training and evaluation | Model conversion (to appropriate format) | Evaluation | Model deployment | Monitoring and observability |

**Code**

Training code

Test code

Application code

**Model**

Candidate models

Chosen model

Offline model

Offline model

Code and model in production

**Data**

Raw data

Labelled data

Training data

Test data    Metrics

Test data

Production data

10

# Implementation of end-to-end lifecycle in AI projects (Alake, 2020), (Sato et al., 2019)

- Problem statement
- Ideal problem solution
- Understanding and insight into the problem
- Technical requirements

- Data structure and source
- Solution form
- Model architecture
- Algorithm research
- Hardware requirements

- Data gathering (diverse, unbiased and abundant)

- Data reformatting
- Data cleaning
- Data normalisation
- Data augmentation

- Usage of pre-trained models?
- Fine-tuning pre-trained models

- Training accuracy
- Validation accuracy
- Training loss
- Validation loss
- Underfitting or overfitting?

- Refine and optimise the model
- Model conversion
- Mobile-optimised model

- Confusion matrix (error matrix)
- Precision-recall

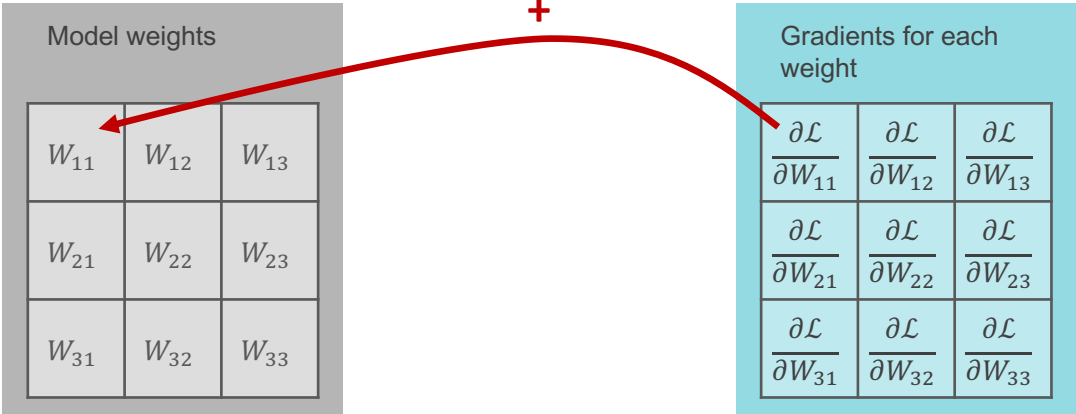- UI interface to access model functionalities
- Continuous integration pipeline that enables model redeployment

- Model performance monitoring system

| Problem definition | Research | Data aggregation, mining and scraping | Data preparation, pre-processing and augmentation | Model building, implementation and experimentation | Model training and evaluation | Model conversion (to appropriate format) | Evaluation | Model deployment | Monitoring and observability |

**Code**

Training code · Test code · Application code

**Model**

Candidate models · Chosen model · Offline model · Offline model · Code and model in production

**Data**

Raw data · Labelled data · Training data · Test data · Metrics · Test data · Production data

# GPT assistant training pipeline (Karpathy, 2023)

- 1000s of GPUs
- Months of training
- Example: GPT, LLaMA, PaLM

→ We can deploy this model.

- 1-100 GPUs
- Days of training
- Example: Vicuna-13B

→ We can deploy this model.

- 1-100 GPUs
- Days of training

- 1-100 GPUs
- Days of training
- Example: ChatGPT, Claude

→ We can deploy this model.

- Refine and optimise the model
- Model conversion
- Mobile-optimised model

- Confusion matrix (error matrix)
- Precision-recall

**Pretraining** | **Supervised Finetuning** | **Reward Modelling** | **Reinforcement Learning** | Model conversion (to appropriate format) | Evaluation

### Code

**Language modelling**
- Predict the next token

**Language modelling**
- Predict the next token

**Binary classification**
- Predict rewards consistent with preferences

**Reinforcement Learning**
- Generate tokens that maximise the reward

Test code

### Model

Base model — Initialised from BM → SFT model — Initialised from SFT → RM model — Initialised from SFT and use RM → RL model → Chosen model → Offline model → Offline model

### Data

**Raw Internet**
- Text trillions of words
- Low-quality and large quantity

**Demonstrations**
- Ideal assistant responses
- ~10K-100K (prompts and responses)
- Written by contractors
- Low quantity and high quality

**Comparisons**
- 100K -1M comparisons
- Written by contractors
- Low quantity and high quality

**Prompts**
- ~10K-100K prompts
- Written by contractors
- Low quantity and high quality

Test data

# GPT assistant training pipeline (Karpathy, 2023)

- 1000s of GPUs
- Months of training
- Example: GPT, LLaMA, PaLM

→ We can deploy this model.

- 1-100 GPUs
- Days of training
- Example: Vicuna-13B

→ We can deploy this model.

- 1-100 GPUs
- Days of training

- 1-100 GPUs
- Days of training
- Example: ChatGPT, Claude

→ We can deploy this model.

- Refine and optimise the model
- Model conversion
- Mobile-optimised model

- Confusion matrix (error matrix)
- Precision-recall

| Pretraining | Supervised Finetuning | Reward Modelling | Reinforcement Learning | Model conversion (to appropriate format) | Evaluation |

**Code**

**Language modelling**
- Predict the next token

**Language modelling**
- Predict the next token

**Binary classification**
- Predict rewards consistent with preferences

**Reinforcement Learning**
- Generate tokens that maximise the reward

Test code

**Model**

Base model  →  Initialised from BM  →  SFT model  →  Initialised from SFT  →  RM model  →  Initialised from SFT and use RM  →  RL model  →  Chosen model  →  Offline model  →  Offline model

**Data**

**Raw Internet**
- Text trillions of words
- Low-quality and large quantity

**Demonstrations**
- Ideal assistant responses
- ~10-100K (prompts and responses)
- Written by contractors
- Low quantity and high quality

**Comparisons**
- 100K -1M comparisons
- Written by contractors
- Low quantity and high quality

**Prompts**
- ~10K-100K prompts
- Written by contractors
- Low quantity and high quality

Test data

# LLMs model sizes over time (Information is Beautiful, 2024)



**Model size (in billion of parameters)**

Labels visible in chart: Gemini*, GPT-5*, Claude-Next*, Wu Dao 2.0, GLaM, BingChat*, Ernie Olympus*, BERT-480, NLG, PaLM, Minerva, PaLM2, Gopher, Titan, Exaone, Titan, GPT-3, PanGu-Alpha, Jurassic-1, Ernie Bot, OPT-IML, BlenderBot3, WebGPT, Tongyi, Ernie, Falcon 180, Q, LaMDA, FLAN, GLM-130B, Galactica, YaLM 100B, Chinchilla, Sparrow, NLLB-200, Vicuna-13B, LLaMA, IDEFICS, Claude, Bloomberg-GPT, Retro48B, Falcon LLM, Grok*, xlarge, Hugging Chat, GPT-NeoX, CodeGen, AlexaTM, PaLI, InternLM, mGPT, Atlas, ChatGPT, Dolly 2, Bio, AceGPT, T5, Megatron 11B, PLATO-XL, BlenderBot, MoLM, mon-8B, Alpha 7B, BlenderBot1, InCoder, FIM 6.9B, GPT4All, Web, DeciLM, GPT-J, Bolt 2.5B

**Trendline**

The model size has exponentially increased, notably achieving approximately 1.8 trillion parameters with the introduction of GPT-4.

# Full-parameter fine-tuning

- Updates are applied to all model weights.
- Models feature large weight matrices, e.g., 7 billion weights for a 7B model and 13 billion for a 13B model.
- Weight updates occur over multiple epochs.
- Extensive memory is required to store and update weights.
- Fine-tuning is restricted to high-capacity GPUs or GPU clusters due to these memory demands.

Model weights

| $W_{11}$ | $W_{12}$ | $W_{13}$ |
|---|---|---|
| $W_{21}$ | $W_{22}$ | $W_{23}$ |
| $W_{31}$ | $W_{32}$ | $W_{33}$ |

**+**

Gradients for each weight

| $\frac{\partial \mathcal{L}}{\partial W_{11}}$ | $\frac{\partial \mathcal{L}}{\partial W_{12}}$ | $\frac{\partial \mathcal{L}}{\partial W_{13}}$ |
|---|---|---|
| $\frac{\partial \mathcal{L}}{\partial W_{21}}$ | $\frac{\partial \mathcal{L}}{\partial W_{22}}$ | $\frac{\partial \mathcal{L}}{\partial W_{23}}$ |
| $\frac{\partial \mathcal{L}}{\partial W_{31}}$ | $\frac{\partial \mathcal{L}}{\partial W_{32}}$ | $\frac{\partial \mathcal{L}}{\partial W_{33}}$ |

Suppose hardware constraints limit our ability to test diverse strategies for enhancing the base model. In that case, Low-Rank Adaption (LoRA) offers two principal methods for solving this problem and can fine-tune LLMs at only a fraction of the cost.

# How is Low-Rank Adaption (LoRA) different? (Hu et al., 2021)

1. We **monitor weight changes** instead of directly updating them.

2. These weight changes are tracked in **two distinct and smaller matrices**, which are multiplied to create a product identical in size to the model's weight matrix.

Original model weights          LoRA weight changes                    Fine-tuned model weights

+          =

# How is Low-Rank Adaption (LoRA) different? (Hu et al., 2021)

1. We **monitor weight changes** instead of directly updating them.

2. These weight changes are tracked in **two distinct and smaller matrices**, which are multiplied to create a product identical in size to the model's weight matrix.

Original model weights                 LoRA low-rank matrices                 Fine-tuned model weights

+   (  [ ]  X  [ ]  )   =

# How is Low-Rank Adaption (LoRA) different? (Hu et al., 2021)

1. We **monitor weight changes** instead of directly updating them.

2. These weight changes are tracked in **two distinct and smaller matrices**, which are multiplied to create a product identical in size to the model's weight matrix.



Original model weights

LoRA matrices, rank 2

x

=

Higher precision weight changes

The precision of the fine-tuning process can be enhanced by increasing the rank.

# Number of trainable parameters

| Rank | Model size (in billion of parameters) | | | |
|------|------|------|------|------|
| | 7B | 13B | 70B | 180B |
| 1 | 167k | 228k | 529k | 849k |
| 2 | 334k | 456k | 1M | 2M |
| 8 | 1M | 2M | 4M | 7M |
| 16 | 3M | 4M | 8M | 14M |
| 512 | 86M | 117M | 270M | 434M |
| 1,024 | 171M | 233M | 542M | 869M |
| 8,192 | 1.4B | 1.8B | 4.3B | 7.0B |

*In reality, LLMs consist of multiple layers of varying sizes, contrary to the simplification of being a single layer.*

# Percent of total parameters

| Rank | Model size (in billion of parameters) | | | |
|---|---|---|---|---|
| | **7B** | **13B** | **70B** | **180B** |
| **1** | 0.00% | 0.00% | 0.00% | 0.00% |
| **2** | 0.01% | 0.00% | 0.00% | 0.00% |
| **8** | 0.02% | 0.01% | 0.01% | 0.00% |
| **16** | 0.04% | 0.03% | 0.01% | 0.01% |
| **512** | 1.22% | 0.90% | 0.39% | 0.24% |
| **1,024** | 2.45% | 1.80% | 0.77% | 0.48% |
| **8,192** | 19.58% | 14.37% | 6.19% | 3.86% |

**Involvement from higher-ranked individuals is particularly beneficial in teaching complex behaviours and addressing behaviours that contradict or fall outside the range of initial training.**

*Percentages may be understated due to the multi-layered structure of models, but the core concept remains clear.*

# QLoRA <span>(Dettmers et al., 2023)</span>
## — Efficient fine-tuning of quantised LLMs

- This is basically LoRA 2.0 with "recoverable" quantisation for reduced memory usage.

- The scientific paper has two critical findings:
  - Training all network layers is crucial for matching the performance of full-parameter fine-tuning.
  - The rank values between 8 and 256 are observed to have minimal impact on performance.

# Overcoming challenges of LLMs

- **Hallucination problem:**
  - Models frequently generate false information with high confidence, presenting a significant risk in scenarios that require strict accuracy. This represents the most prominent challenge in Generative AI.

- **Attribution problem:**
  - More clarity is needed regarding why models produce specific outputs, which makes it difficult to trust or validate their responses.

- **Staleness:**
  - Language models quickly become outdated, needing more information on recent events, diminishing their relevance and utility over time.

- **Revisions challenge:**
  - Models must comply with regulations (e.g., GDPR), which entails the ability to delete or revise data, a functionality that remains underdeveloped. The AI Act commits to monitoring across various dimensions of risk: fairness, autonomy, transparency, security, reliability, and data protection.

- **Customisation issue:**
  - Adapting models to specific use cases or datasets is an on-going challenge, necessitating innovative solutions for effective personalisation and application in diverse environments. This includes strategies for integrating models with unique corporate data or adjusting outputs to align with the specific needs of different user groups.

A prevalent strategy currently being adopted involves integrating existing LLMs with external memory resources. Retrieval-Augmented Generation (RAG) systems, which dynamically retrieve and incorporate external data into the decision-making process, are solving these challenges.

# General Purpose AI (GPAI) classification and key requirements for providers (European Parliament, 2024), (Pinto, 2024)

## ALL GPAI MODELS

Large models and systems capable of competently performing a wide range of distinctive tasks, such as generating video, text, images, computer code, or conversing.

- Transparency obligations before market placement, including:
  - Drawing up technical documentation for downstream providers
  - Complying with EU copyright law and disseminating detailed summaries about the content used in training
  - Watermarking AI generated or manipulated content

## SYSTEMIC RISK GPAI MODELS

Foundation models trained with a large amount of data and with advanced complexity, capabilities, and performance well above the average can disseminate systemic risks along the value chain.

- Complying with all requirements applicable to all GPAI models and systems
- Conducting model evaluations
- Assessing and mitigating systemic risks
- Conducting adversarial testing
- Reporting of serious incidents to the EU Commissions
- Ensuring sufficient cybersecurity protection
- Reporting on energy efficiency

# Fine Tuning vs. Retrieval-Augmented Generation (RAG) (Soudani et al., 2024)

- Researchers have studied the comparison between Retrieval Augmented Generation (RAG) and fine-tuning methods on synthetic data.
  - Their investigations reveal that both strategies significantly enhance the capability of AI to handle specialised information during question-answering tasks.
  - **RAG emerges as the leading methodology, outperforming fine-tuning in improving model responses to obscure queries.**
  - This does not eliminate fine-tuning's relevance but suggests **RAG as a more efficient option for bolstering AI against niche topics**.
  - Fine-tuning is acknowledged for its depth in embedding knowledge, but it shares similar limitations with pre-training, particularly in learning about infrequent concepts.

# Implementation of end-to-end lifecycle in AI projects (Alake, 2020), (Sato et al., 2019)

- Problem statement
- Ideal problem solution
- Understanding and insight into the problem
- Technical requirements

- Data structure and source
- Solution form
- Model architecture
- Algorithm research
- Hardware requirements

- Data gathering (diverse, unbiased and abundant)

- Data reformatting
- Data cleaning
- Data normalisation
- Data augmentation

- Usage of pre-trained models?
- Fine-tuning pre-trained models

- Training accuracy
- Validation accuracy
- Training loss
- Validation loss
- Underfitting or overfitting?

- Refine and optimise the model
- Model conversion
- Mobile-optimised model

- Confusion matrix (error matrix)
- Precision-recall

- UI interface to access model functionalities
- Continuous integration pipeline that enables model redeployment

- Model performance monitoring system

| Problem definition | Research | Data aggregation, mining and scraping | Data preparation, pre-processing and augmentation | Model building, implementation and experimentation | Model training and evaluation | Model conversion (to appropriate format) | Evaluation | Model deployment | Monitoring and observability |

**Code**

Training code · Test code · Application code

**Model**

Candidate models · Chosen model · Offline model · Offline model · Code and model in production

**Data**

Raw data · Labelled data · Training data · Test data · Metrics · Test data · Production data

25

# Basic chatbot architecture

# Basic chatbot architecture (Simon, 2023)
— Example



Model deployment

User

Code

Model

Data

"What is the latest trend for solar investments in China?"

Application code

"As a helpful energy specialist, please answer the question, focusing on numerical data. Do not invent facts. If you cannot provide a factual answer, say you do not know the answer."

"According to a report by the International Energy Agency (IEA), China was the world's largest solar market in 2020, with a total installed capacity of 160 GW. The report also states that China's solar market is expected to continue to grow, with a target of 250 GW of installed capacity by 2025. However, the report does not provide specific information on the latest trend for solar investment in China."

Document corpus

27

# Retrieval-Augmented Generation (RAG) architecture

# Retrieval-Augmented Generation (RAG) architecture (Simon, 2023)
— Example



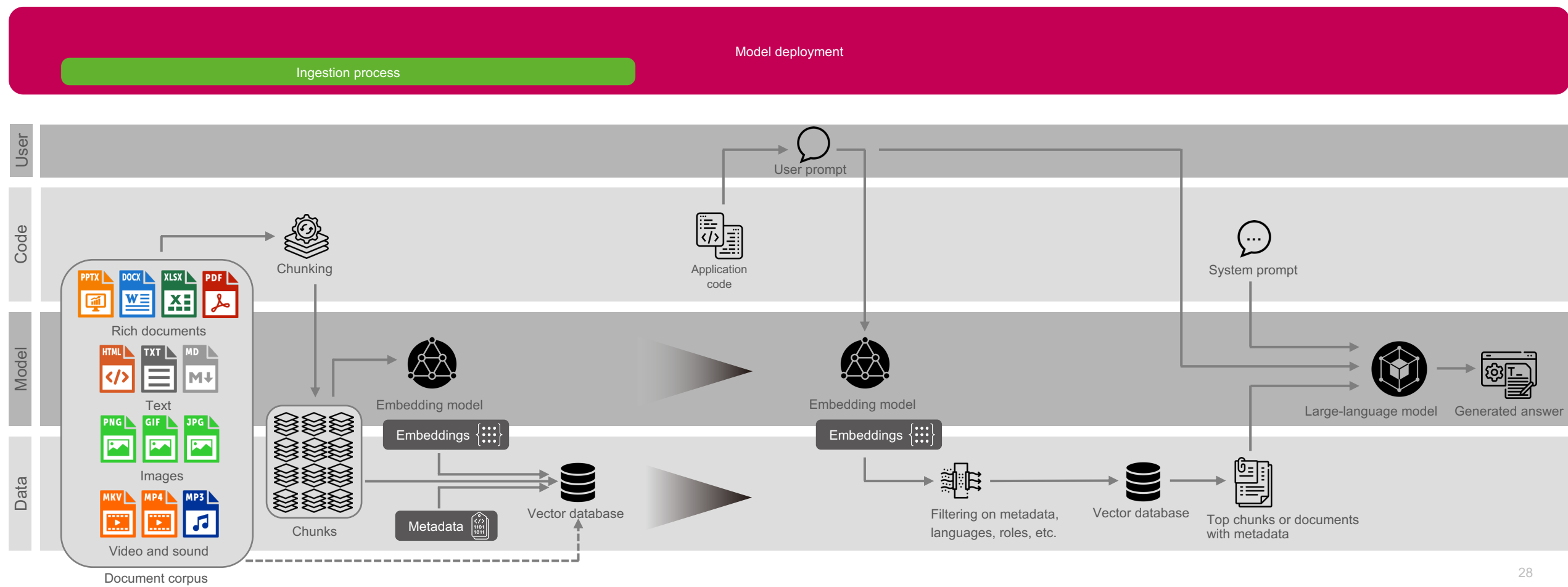**Model deployment**

**Ingestion process**

"What is the latest trend for solar investments in China?"

Application code

"As a helpful energy specialist, please answer the question, focusing on numerical data. Do not invent facts. If you cannot provide a factual answer, say you do not know the answer.

Question: {question}

Useful context to expand your build-in knowledge: {context}"

- WorldEnergyInvestment2023.pdf
- WorldEnergyOutlook2023.pdf

Chunking

Text

Images

Video and sound

Document corpus

Embedding model

Embeddings

Metadata

Chunks

Vector database

Embedding model

Embeddings

Filtering on languages,

"The latest trend for solar investment in China is that solar PV capacity additions in China reached over 270 GW per year before flattening in the STEPS scenario. This represents a marked slowing of the rate of growth achieved from 2023 to 2050 in the NZE scenario. China remains the largest solar PV market, accounting for 45% of all capacity additions in 2022".

User
Code
Model
Data

29

# Retrieval-Augmented Generation (RAG) architecture (Simon, 2023)

— Example



Model deployment

Ingestion process

User

"What is the latest trend for solar investments in China?"

"What does STEPS mean"

"As a helpful energy specialist, please answer the question, focusing on numerical data. Do not invent facts. If you cannot provide a factual answer, say you do not know the answer.
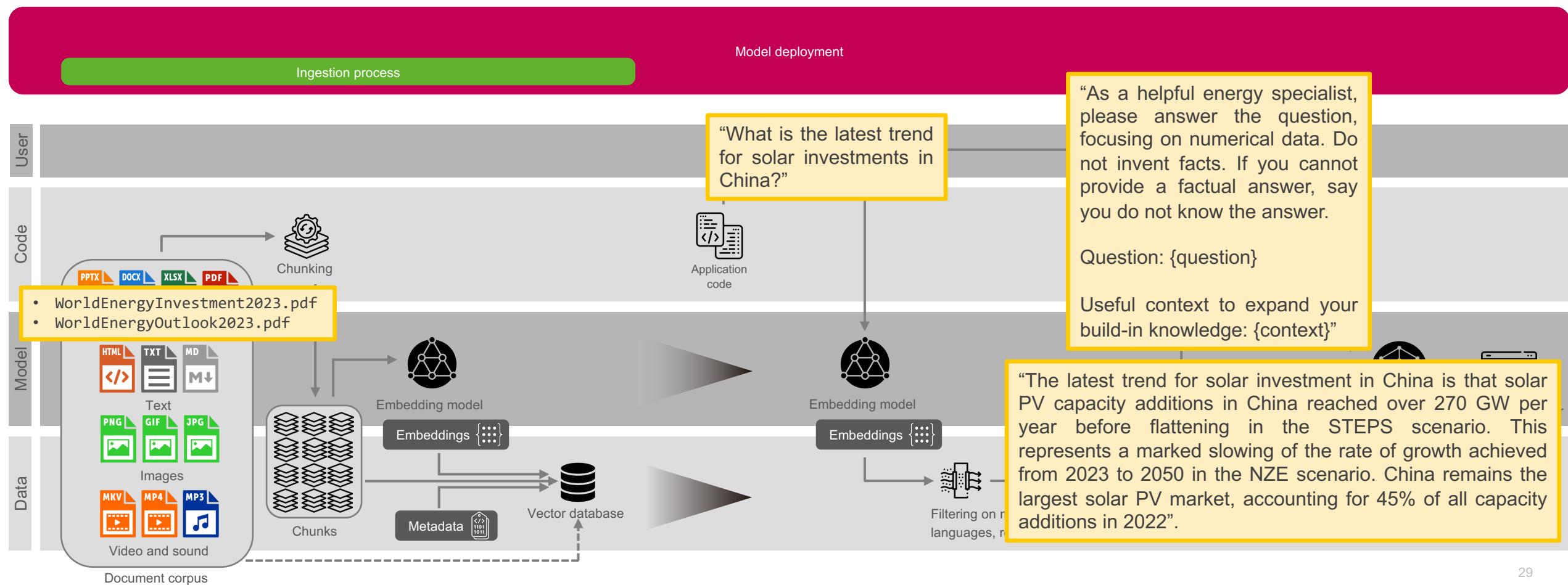
Question: {question}"

Code

Application code

Chunking

- WorldEnergyInvestment2023.pdf
- WorldEnergyOutlook2023.pdf

PPTX  DOCX  XLSX  PDF

Model

HTML  TXT  MD

Text

Embedding model

Embedding model

Embeddings

Embeddings

"The STEPS scenario is a scenario that provides a sense of the prevailing direction of energy system progression, based on a detailed review of the current policy landscape. It explores how energy systems evolve under current conditions and gives a sense of the current direction of travel for the energy economy. The STEPS is associated with a sector-by-sector review of the policies and measures that are actually in place or that have been announced; aspirational energy or climate targets are not automatically assumed to be met. ..."

Data

PNG  GIF  JPG

Images

MKV  MP4  MP3

Video and sound

Chunks

Metadata

Vector database

Filtering on languages, r

additions in 2022".

Document corpus

30

# Retrieval-Augmented Generation (RAG) architecture
— Many questions ?!?

**How to scale?**

**How to learn?**

**How to optimise?**

Model deployment

Ingestion process

**How to encode queries?**

**How to encode?**

**How to chunk?**

**How to prompt?**

**How to post-process?**

User prompt

**How to pass context?**

**How to verify?**

Application code

System prompt

**When to retrieve?**

Chunking

PPTX DOCX XLSX PDF
Rich documents

HTML TXT MD
Text

Embedding model

Embeddings

Embedding model

Embeddings

**How and what to retrieve?**

Large-language model

Generated answer

PNG GIF JPG
Images

MKV MP4 MP3
Video and sound

Document corpus

Chunks

Metadata

Vector database

Filtering on metadata, languages, roles, etc.

Vector database

Top chunks or documents with metadata

**How to pre-process?**

User

Code

Model
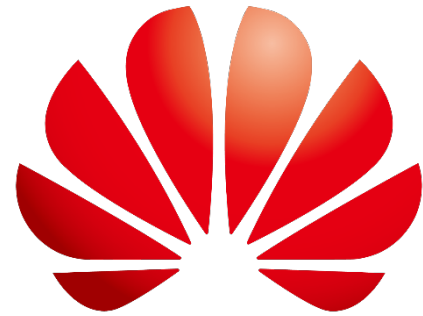
Data

31

# Retrieval-Augmented Generation (RAG) architecture

**Chunking strategy:**
- Chunk size optimisation
- Character, recursive character or document specific
- Sliding window chunking

**Embedding strategy:**
- Which embedding algorithm or model?

(Muennighoff et al., 2023), (Hugging Face, 2024)

**Document retriever:**
- Metadata attachment
- Mixed retrieval
- Cognitive reviewer

**Failure points:**
- Missing content
- Missed the top-ranked documents
- Not in context – consolidation strategy limitations
- Not extracted
- Wrong format

**User authentication:**
- Access control
- Data security
- User privacy
- Legal compliance
- Accountability

**Guardrails:**
- Anonymization
- Restrict substrings, topics, code, language
- Detect prompt injection
- Detect toxicity

**Query strategy:**
- Rewrite based on history
- Create subqueries or similar queries
- Query cost

**Choice of LLM:**
- Architecture
- Number of parameters
- Access
- Use-case
- Data privacy

**Evaluate responses:**
- Prompt evaluation
- RAG retrieval evaluation
- Relevance metrics
- Tasks-specific metrics
- Alignment metrics

Model deployment

Ingestion process

**User**

User prompt

**Code**

Application code

System prompt

Chunking

**Model**

Rich documents

PPTX DOCX XLSX PDF

HTML TXT MD

Text

Embedding model

Embeddings

Embedding model

Embeddings

Large-language model

Generated answer

**Data**

PNG GIF JPG

Images

MKV MP4 MP3

Video and sound

Document corpus

Chunks

Metadata

Vector database

Filtering on metadata, languages, roles, etc.

Vector database

Top chunks or documents with metadata

# Closing remarks

- A single 7-day forecast consumes 14 Wh with Huawei Pangu-Weather compared to 30,000 Wh with the ICON model, illustrating a significant difference in energy efficiency.

- LoRA fine-tunes LLMs by monitoring and updating weight changes through smaller matrices, enhancing fine-tuning precision without direct weight modification.

- LLMs can generate false information with high confidence, presenting a significant risk in scenarios that require strict accuracy.

- Integrating existing LLMs with external memories through Retrieval-augmented Generation (RAG) systems is a leading solution to current challenges.

- RAG outshines fine-tuning in AI's handling of niche topics by significantly enhancing response precision to obscure queries.

- Despite RAG's superiority in handling niche topics, ongoing research and numerous open questions highlight the evolving nature of this AI methodology.

Advancing the Intelligent World

# References

Alake, R. (2020). *10 Stages Of A Machine Learning Project In 2020 (And Where You Fit)*. *2020*, 1–17.

Dettmers, T., Pagnoni, A., Holtzman, A., & Zettlemoyer, L. (2023). *QLoRA: Efficient Finetuning of Quantized LLMs* (arXiv:2305.14314). arXiv. https://doi.org/10.48550/arXiv.2305.14314

European Parliament. (2024). *Provisional Agreement Resulting from Interinstitutional Negotiations*. https://doi.org/10.5040/9781782258674

Information is Beautiful. (2024). *The Rise of Generative AI Large Language Models (LLMs) like ChatGPT*. Information Is Beautiful. https://informationisbeautiful.net/visualizations/the-rise-of-generative-ai-large-language-models-llms-like-chatgpt/

Karpathy, A. (2023, May 23). *State of GPT*. Microsoft BUILD. https://karpathy.ai/stateofgpt.pdf

Pinto, T. (2024, January 8). General-purpose AI under the AI Act. *Artificial Intelligence Act*. https://artificialintelligenceact.com/general-purpose-ai-under-the-ai-act/

Sato, D., Wider, A., & Windheuser, C. (2019). *Continuous Delivery for Machine Learning*. https://martinfowler.com/articles/cd4ml.html

Simon, J. (Director). (2023, October 24). *Retrieval-Augmented Generation chatbot, part 1: LangChain, Hugging Face, FAISS, AWS*. https://www.youtube.com/watch?v=7kDaMz3Xnkw