# Agenda

Evaluation and Benchmark

Parametric Knowledge Adaptation ~60min
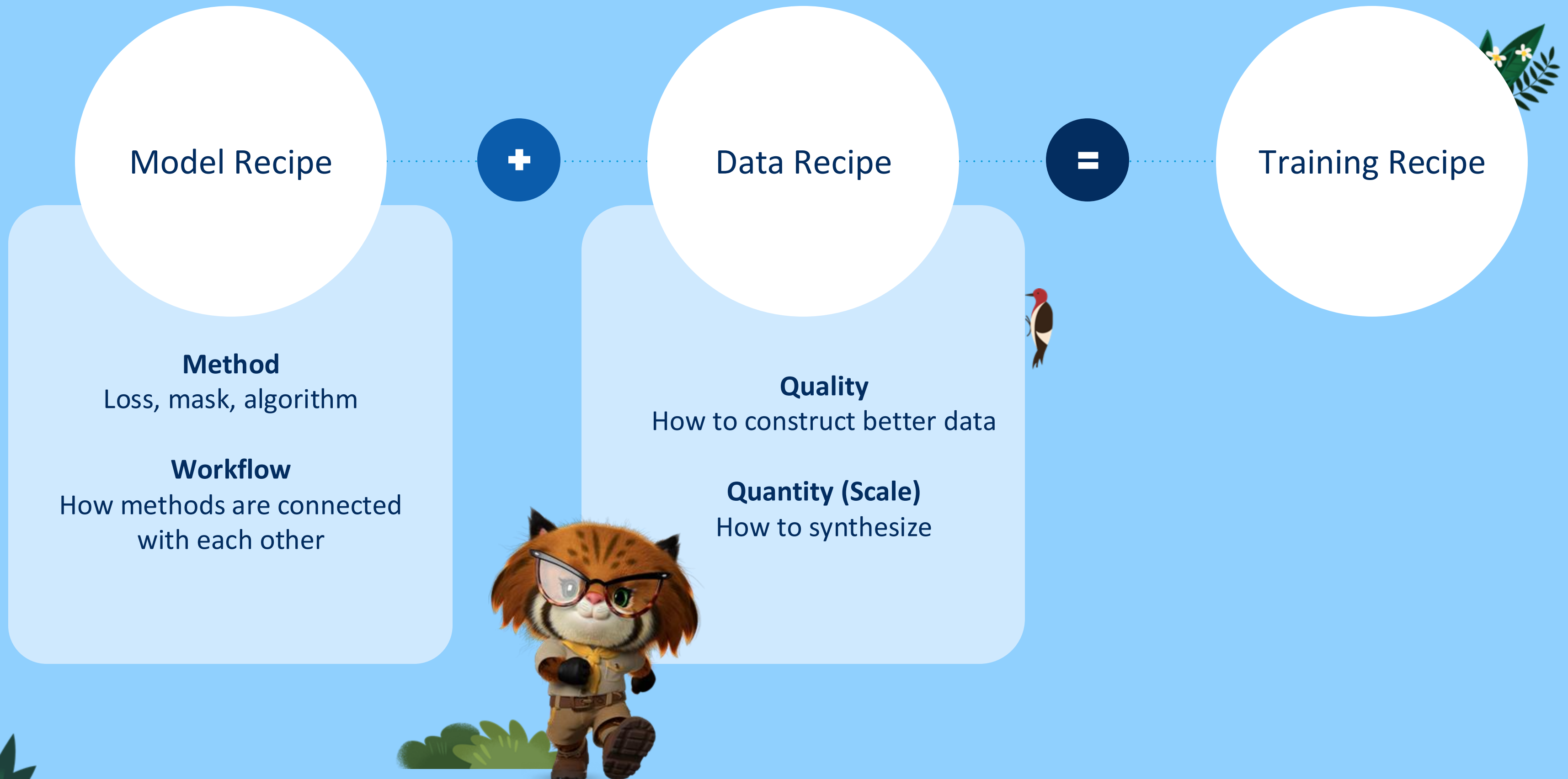
Semi-Parametric Knowledge Adaptation

Summary, Discussion, QAs

# Adaptation - Overview

**Model Recipe**  +  **Data Recipe**  =  **Training Recipe**

**Method**
Loss, mask, algorithm

**Workflow**
How methods are connected with each other

**Quality**
How to construct better data

**Quantity (Scale)**
How to synthesize

# Adaptation - Overview

## Training Recipe

**Data Recipe:**
 e.g., Supervised data is expensive, how to synthesize more data?

**Model Recipe:**
 e.g., **Hyper-parameters**: What are the important hyper-parameters?

 e.g., **Training Workflow**: How to connect with other methods?

## Seed Data

**Data Acquisition:**
e.g., crawling, quality, quantity, filtering…

**Data Mixture:**
e.g., in-domain, general-domain, …

**Data Budget:**
e.g., instruction following ~ 1 million; preference learning ~ 1 million (often overlapping with instruction following prompt); reinforcement learning ~ 10-100 thousand

Ke, Ming, Joty - Adaptation of LLMs Tutorial, NAACL 2025

# Continual Pre-training (CPT)

# CPT – Role

## Knowledge Transfer

**Improves on new knowledge:**

CPT is typically used to inject new knowledge/capability (e.g., long-context adaptation) to the base model and to provide good initialization to the subsequent stages
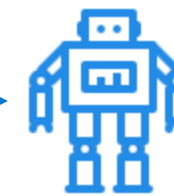
## Prevent Forgetting

**Reinforce similar problems:**

CPT involves large amount of unsupervised data and could easily cause *catastrophic forgetting* to the base model

Ke, Ming, Joty - Adaptation of LLMs Tutorial, NAACL 2025

# CPT – Example Workflow

Seed Data (unsupervised)



Next Token Prediction*
(self-supervised)

*Potentially some modifications (e.g., position embedding modification in long-context adaptation)

Ke, Ming, Joty - Adaptation of LLMs Tutorial, NAACL 2025

6

# CPT – Example Data

Long Text
(e.g. website, books)

No Special Masking

# CPT – Key Considerations

## Training Recipe

**Model Recipe:**
   **Hyper-parameters**: What are the important hyper-parameters?

   **Training Workflow**: how to connect CPT with other methods (e.g., IT, SPL)

## Seed Data

**Data Source:** Where to get the data?

**Data Mixture:** What should be included to the CPT data?
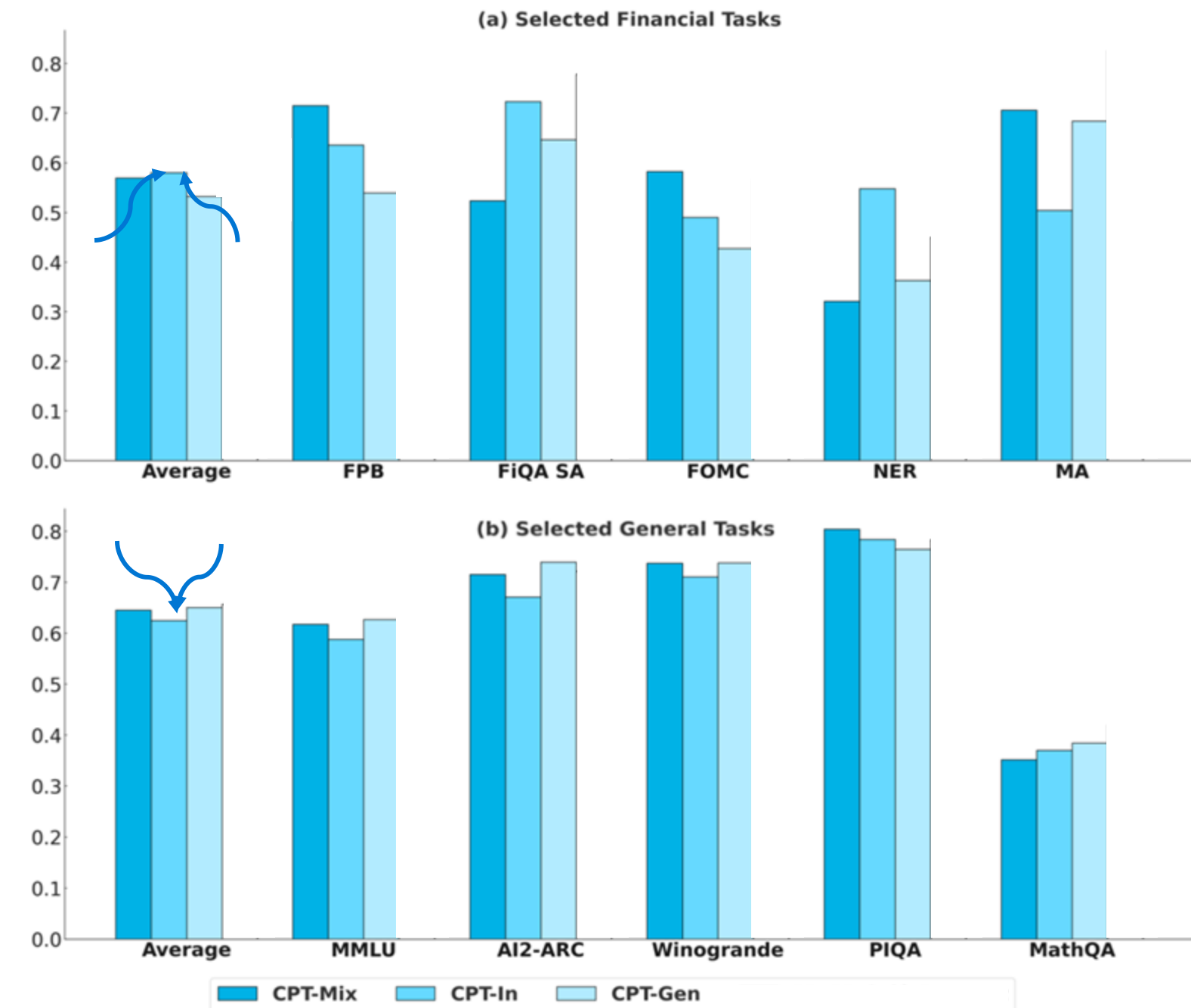
**Data Budget:** How much data we need?

Ke, Ming, Joty - Adaptation of LLMs Tutorial, NAACL 2025

# CPT – Key Ideas

## Catastrophic Forgetting (Finance-LLM as an example)

In-domain Data alone → forgetting on
general knowledge
(Knowledge forgetting)

Demystifying Domain-adaptive Post-training for Financial LLMs, Ke et al., 2025



(a) Selected Financial Tasks

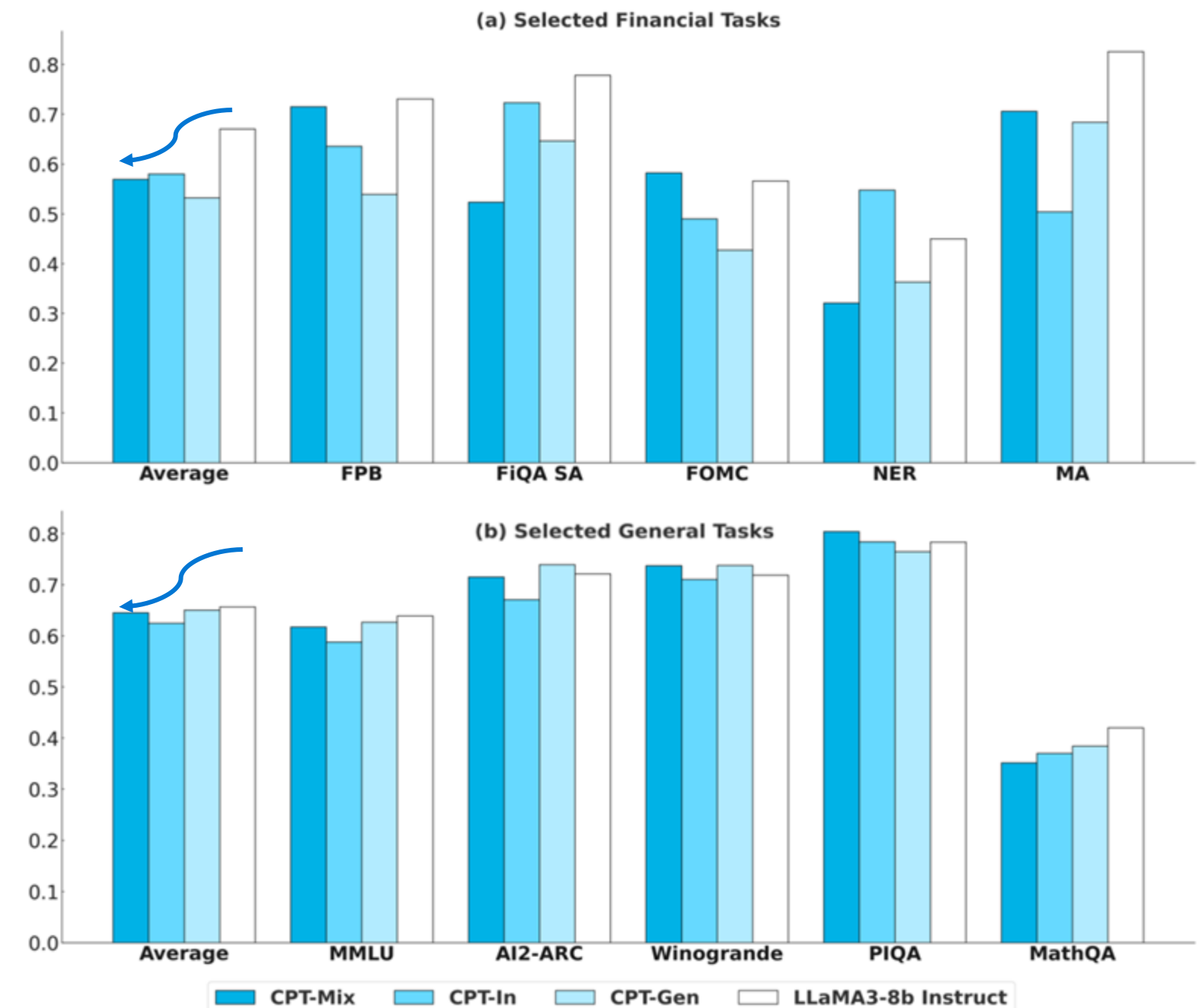(b) Selected General Tasks

CPT-Mix    CPT-In    CPT-Gen

# CPT – Key Ideas

## Catastrophic Forgetting (Finance-LLM as an example)

CPT alone →
forgetting on general capabilities
(Capabilities forgetting)

base model = instruction-tuned model

Demystifying Domain-adaptive Post-training for Financial LLMs, Ke et al., 2025



(a) Selected Financial Tasks

(b) Selected General Tasks

CPT-Mix   CPT-In   CPT-Gen   LLaMA3-8b Instruct

Ke, Ming, Joty - Adaptation of LLMs Tutorial, NAACL 2025

10

# "We find that even small amounts of replay (1% of the general domain data) mitigate forgetting

**Demystifying Domain-adaptive Post-training for Financial LLMs**

Zixuan Ke, Yifei Ming, Xuan-Phi Nguy
Salesforce AI
{zixuan.ke,yifei.ming,xnguyen,c

Project Page: https://github.com

Datasets: https://huggingface.cc

**Simple and Scalable Strategies to Continually Pre-train Large Language Models**

Adam Ibrahim[*†©]
Benjamin Thérien[*†©]
Kshitij Gupta[*†©]
Mats L. Richter[†©]
Quentin Anthony[◊†©]
Timothée Lesort[†©]
Eugene Belilovsky[‡©]
Irina Rish[†©]

**Fine-tuned Language Models are Continual Learners**

Thomas Scialom[1*]    Tuhin Chakrabarty[2*]    Smaranda Muresan[2]

[1]Meta AI
[2]Department of Computer Science, Columbia University

tscialom@fb.com, tuhin.chakr@cs.columbia.edu, smara@cs.columbia.edu

# CPT – Key Ideas

## Learn New Knowledge and Mitigate Knowledge Forgetting – Data

**Data source for new domain:**

**Web scrapers** (often the largest proportion of data): e.g., Internet

**User-provided content** (often smaller proportion, but higher-quality): e.g.,. Wikipedia, arXiv,

**Open Publishers** (often smaller proportion, but higher-quality)**:** e.g., PubMed, Semantic Scholar, Text book

**Data source to prevent forgetting (small amount of replay):**

**Human Verifier Text** (small but high-quality)**:** e.g., general supervised tasks

# CPT – Key Ideas

## Learn New knowledge and Mitigate Knowledge Forgetting – Data

General Domain data
\+  In-domain data

| Capability | Domain | CPT Dataset | Size | Reference |
|---|---|---|---|---|
| Concept | General | NaturalInstrution | 100,000 | Mishra et al. (2022) |
| | | PromptSource | 100,000 | Bach et al. (2022) |
| | | Math | 29,837 | Amini et al. (2019b) |
| | | Aqua | 97,500 | Ling et al. (2017) |
| | | CREAK | 10,200 | Onoe et al. (2021) |
| | | ESNLI | 549,367 | Camburu et al. (2018) |
| | | QASC | 8,130 | Khot et al. (2020) |
| | | SODA | 1,190,000 | Kim et al. (2022) |
| | | StrategyQA | 2,290 | Geva et al. (2021) |
| | | UnifiedSKG | 779,000 | Xie et al. (2022) |
| | | GSM8K | 7,470 | Cobbe et al. (2021) |
| | | ApexInstr | 1,470,000 | Huang et al. (2024b) |
| | | DeepmindMath | 379,000 | Saxton et al. (2019) |
| | | DialogueStudio | 1,070,000 | Zhang et al. (2023) |
| | Finance | Fineweb-Fin | 4,380,000 | - |
| | | Book-Fin | 4,500 | - |
| Total | | | 10,177,294 | |

Demystifying Domain-adaptive Post-training for Financial LLMs, Ke et al., 2025

# CPT – Key Ideas

## Learn New knowledge and Mitigate Capabilities Forgetting – Model

**Replay data only addresses the domain knowledge forgetting, but it does not address the capabilities (e.g., instruction-following abilities)**

One way is to jointly train CPT and IT to avoid the capabilities forgetting

- Mitigate forgetting
- Encourage transfer (concept learned from CPT naturally shared across tasks)

**Demystifying Domain-adaptive Post-training for Financial LLMs**

Zixuan Ke, Yifei Ming, Xuan-Phi Nguyen, Caiming Xiong and Shafiq Joty

Salesforce AI Research

{zixuan.ke,yifei.ming,xnguyen,cxiong,sjoty}@salesforce.com

🧠 Project Page: https://github.com/SalesforceAIResearch/FinDAP

🤗 Datasets: https://huggingface.co/datasets/Salesforce/FinEval

\* Another way could be model merging

A SURVEY ON POST-TRAINING OF LARGE LANGUAGE MODELS, Tie et al., 2025

Ke, Ming, Joty - Adaptation of LLMs Tutorial, NAACL 2025

# CPT – Key Ideas

## Other Tips: Learning Rate, Data Curriculum

### Final Recipe for Llama-Fin

**Continual Pre-training (CPT) and Instruction Tuning (IT)**

| | | |
|---|---|---|
| **Data** | 50% CPT, 50% IT | |
| **Curriculum** | Group 1 | CPT: 50% Domain-specific Text (Web and book), 50% General text (verfiable text)<br>IT: 20% Domain-specific tasks, 80% General tasks |
| | Group 2 | CPT: Group 1 data + domain-specific books<br>IT: Group1 + Exercises extracted from books |
| **Steps** | | Group 1: 3.84B tokens; Group 2: 1.66B tokens<br>(8,000 context length, 16 A100) |
| **Model** | Intialization | Llama3-8b-instruct |
| | Attention | CPT: full attention with cross-docuemnt attention masking<br>IT: full attention with instruction mask-out and cross-docuemnt attention masking |
| **Optim.** | | AdamW (weight decay = 0.1, $\beta_1$=0.9, $\beta_2$=0.95) |
| | LR | Group 1: 5e-6 with 10% warmup; Group 2: 5e-6 with 50% warmup |
| | Batch size | 128K tokens |
| **Stop Cri.** | Loss of development set stops decreasing ($\approx$ 1 epoch) | |

Demystifying Domain-adaptive Post-training for Financial LLMs, Ke et al., 2025

Ke, Ming, Joty - Adaptation of LLMs Tutorial, NAACL 2025

# CPT – Key Ideas

## Other Tips: Learning Rate, Data Curriculum

| Continued Long-context Training | | |
|---|---|---|
| **Data** | 30% code repos, 30% books, 3% textbooks, 37% ShortMix | |
| | ShortMix: | 27% FineWeb-Edu, 27% FineWeb, 11% Wikipedia, 11% StackExchange, 8% Tulu-v2, 8% OpenWebMath, 8% ArXiv |
| **Length Curriculum** | Stage 1 (64K): | Code repos, books, and textbooks at length 64K |
| | Stage 2 (512K): | Code repos: 50% at length 512K, 50% at length 64K Books: 17% at length 512K, 83% at length 64K Textbooks at length 512K |
| **Steps** | Stage 1: 20B tokens (2.2K H100 hours),   Stage 2: 20B tokens (12.2K H100 hours) | |
| **Model** | Initialization: | Llama-3-8B-Instruct (original RoPE base freq. $5 \times 10^5$) |
| | RoPE: | Stage 1: $8 \times 10^6$, Stage 2:   $1.28 \times 10^8$ |
| | Attention: | Full attention with cross-document attention masking |
| **Optim.** | AdamW (weight decay = 0.1, $\beta_1 = 0.9$, $\beta_2 = 0.95$) | |
| | LR: | $1e - 5$ with 10% warmup and cosine decay to $1e - 6$, each stage |
| | Batch size: | 4M tokens for stage 1, 8M tokens for stage 2 |

How to Train Long-Context Language Models (Effectively), Gao et al., 2025

# CPT – Key Ideas

## Other Tips: Learning Rate, Data Curriculum

### Rules of thumb for continual pre-training

**Caveat**—The following guidelines are written to the best of our *current knowledge.*

**Learning rate schedule:**

- If the learning rate was cosine-decayed from a large value $\eta_{max}$ to a small value $\eta_{min}$ during pre-training on the initial dataset, the following guidelines can help to continually pre-train your model:

  - Re-warming and re-decaying the learning rate from $\mathcal{O}(\eta_{max})$ to $\mathcal{O}(\eta_{min})$ improves adaptation to a new dataset, e.g. compared to continuing from small learning rates $\mathcal{O}(\eta_{min})$.
  - Decreasing the schedule's maximum learning rate can help reduce forgetting, whereas increasing it can improve adaptation.

- Infinite LR schedules are promising alternatives to cosine decay schedules. They transition into a high constant learning rate across tasks, helping prevent optimization-related forgetting by avoiding re-warming the LR between tasks. They also avoid committing to a specific budget of tokens as a final exponential decay can be used to train the model to convergence at any point during training.

Simple and Scalable Strategies to Continually Pre-train Large Language Models, Ibrahim et al., 2024

Ke, Ming, Joty - Adaptation of LLMs Tutorial, NAACL 2025

# CPT – Key Ideas

## Other Tips: Learning Rate, Data Curriculum



**Recipe**

- Start with a data distribution that is similar to the pretraining set but places larger weight on high quality sources before transitioning to a second distribution that incorporates QA data and upweights sources in areas of model weakness.

- The learning rate schedule should start from $\eta_{min}$ of the pretrained model and decay with cosine annealing to $\frac{\eta_{min}}{100}$.

- The switch between data distribution should occur at $\frac{\eta_{max}}{5}$ in the learning rate schedule.

Reuse, Don't Retrain: A Recipe for Continued Pretraining of Language Models, Parmar et al., 2024

Ke, Ming, Joty - Adaptation of LLMs Tutorial, NAACL 2025

# CPT – Key Ideas Summary

## Training Recipe

**Model Recipe:**
    **Learning rate schedule**
    **Data curriculum**

    **Jointly training CPT and IT have been shown to be effective**

## Seed Data

**Data Mixture:** Wide representative and filtering is needed

**Data Budget:**
    **New Knowledge ~** 5 million
    **Prevent Forgetting ~** 5 million

\* Filtering can be complicated and involved different components (e.g., decontamination..).

Opening the Language Model Pipeline: A Tutorial on Data Preparation, Model Training, and Adaptation, NeurIPS 2025

# Instruction Tuning

# IT – Role

| Chat Style Adaptation | Chat Template Adaptation |
|---|---|
| Adapt base model to **specific style of input** for chat interactions. | Ability to include **system prompts, multi-turn dialogues,** and other **chat templates.** |

Special tokens

```
<|system|>
You are a helpful assitant
<|end|>
<|user|>
How many helicopters can you eat?
<|end|>
<|assistant|>
{Answer goes here}
```

System prompt

Multi-turn dialogue

Ke, Ming, Joty - Adaptation of LLMs Tutorial, NAACL 2025

# IT – Example Workflow



A SURVEY ON POST-TRAINING OF LARGE LANGUAGE MODELS, Tie et al., 2025

# IT – Example Data

Chat Format
Special Label Masking
Packing



Ke, Ming, Joty - Adaptation of LLMs Tutorial, NAACL 2025

# IT – Key Considerations

## Training Recipe

**Data Recipe:**
Supervised data is expensive, how to synthesize more data?

**Model Recipe:**
How should the loss and masking different from CPT?

**Training Workflow**: how to connect with other methods

## Seed Data

**Data Source:** Where to get the data?

**Data Mixture:** What should be included in the IT data?

**Data Budget:** How many data we need?

Ke, Ming, Joty - Adaptation of LLMs Tutorial, NAACL 2025

## Self-instruct / Synthetic data

**Seed:** N high-quality (often human) prompts

**Ask a strong LLM:** Create a modified version of these instructions

**Generate completions** with another (or same) strong LLM.

**Results:** easily 10x more synthetic training data



Alpaca: A Strong, Replicable Instruction-Following Model, Taori et al., 2023
SELF-INSTRUCT: Aligning Language Models with Self-Generated Instructions, Wang et al., 2022

Ke, Ming, Joty - Adaptation of LLMs Tutorial, NAACL 2025

# IT – Key Ideas

## Packing and Label Masking



https://github.com/MeetKai/functionary/blob/main/functionary/train/packing

# IT – Key Ideas

## Packing and Label Masking

**Disabling cross-document attention.** Ding et al. (2024a) show that masking out attention across document boundaries improve model performance and this was also used during Llama-3 pre-training (Dubey et al., 2024). In §B.2, we show that disabling cross-document attention in continued training benefits both the short and long-context performance. Disabling cross-document attention can also result in higher training throughput, which we describe in more detail in §A.3.

**Packing** Packing optimizes the training efficiency by grouping sequences of varying lengths into a single long sequence without requiring any padding. This technique, commonly used in LLM pre-training, is now also utilized in instruction-based supervised fine-tuning, as implemented by models like Zephyr (Tunstall et al., 2023b)[4].

**Papers show that packing is helpful**

How to Train Long-Context Language Models (Effectively), Gao et al., 2025
LIONs: An Empirically Optimized Approach to Align Language Models, Yu et al., 2024

# IT – Key Ideas

## Packing and Label Masking



**Masking the tokens of the instruction by setting the token labels of the instructions to -100**

https://www.linkedin.com/pulse/llm-research-insights-instruction-masking-new-lora-raschka-phd-7p1oc

Ke, Ming, Joty - Adaptation of LLMs Tutorial, NAACL 2025

# IT – Key Ideas

## Packing and Label Masking



**RQ1: What is the role of DAPT and SFT in post-training?**

- DAPT uses next-token prediction, while SFT needs instruction masking added. §5.1
- Both DAPT and SFT contribute to improvements. §5.2
- Joint training with DAPT and SFT yields better results than sequential training. §5.3

**Papers show that label masking is helpful**

**Loss Masking**   The standard language model training computes loss across all tokens in a sequence. Loss masking, however, ignores loss computation on tokens that are not output tokens like user instructions. It prevents the model from learning irrelevant information, alleviating catastrophic forgetting and overfitting.

Demystifying Domain-adaptive Post-training for Financial LLMs, Ke et al., 2025
LIONs: An Empirically Optimized Approach to Align Language Models, Yu et al., 2024
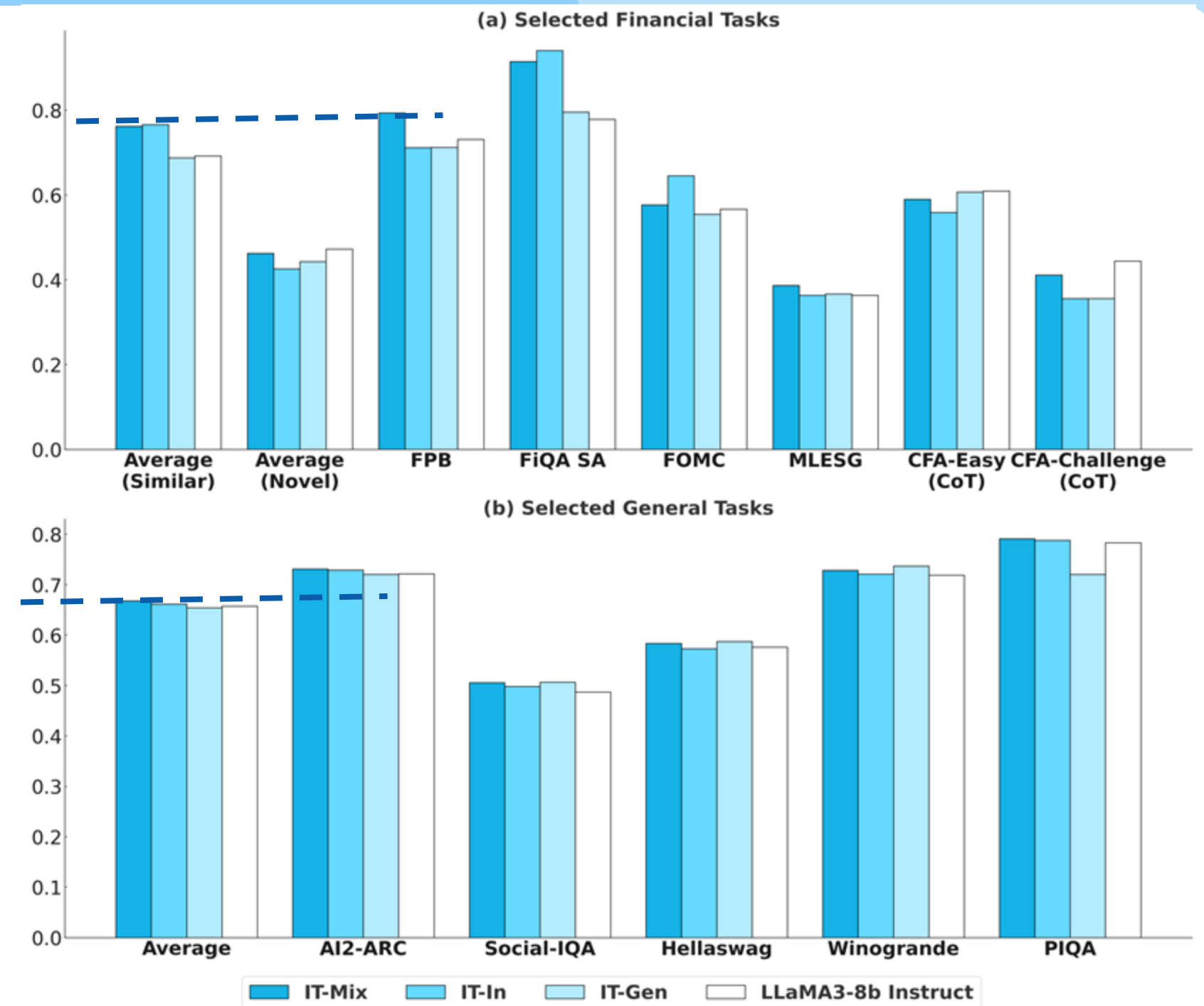
# IT – Key Ideas

## Task Generalization



**Forgetting is less a problem**

**Task generalization is the main issue.**

Demystifying Domain-adaptive Post-training for Financial LLMs, Ke et al., 2025
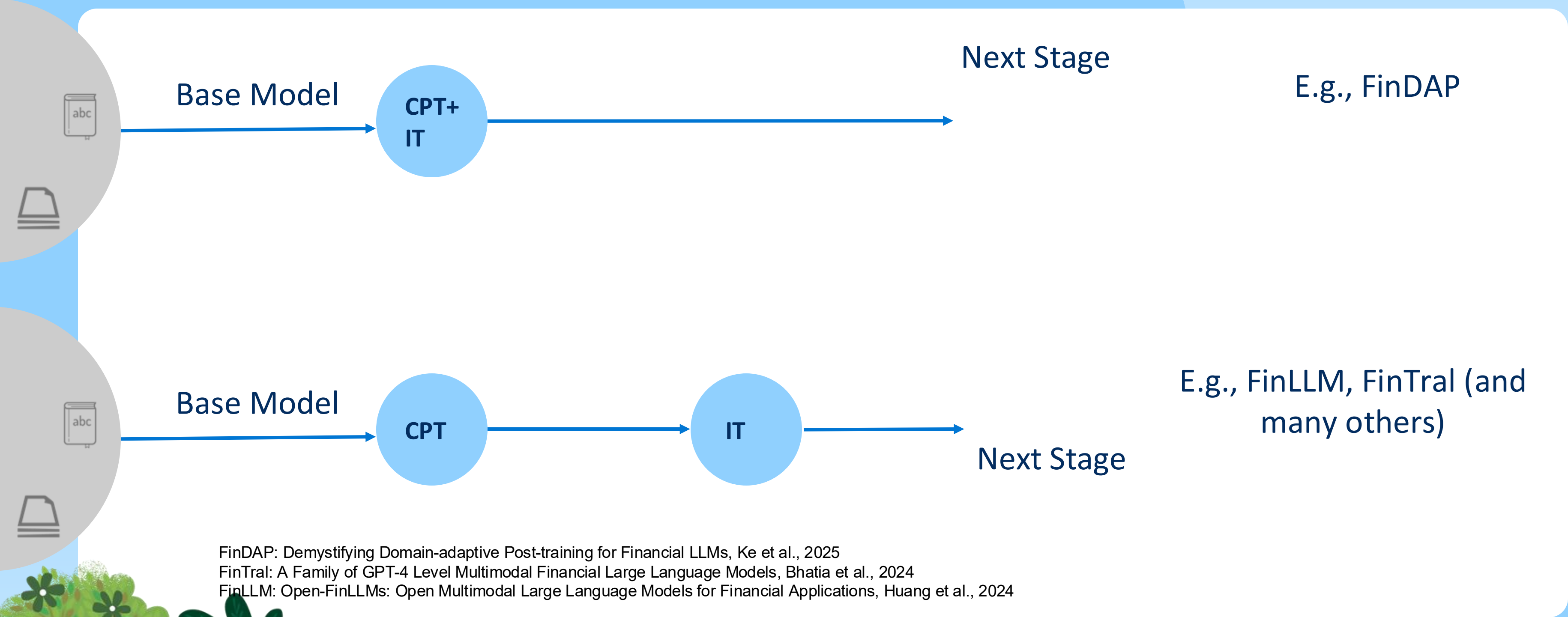
# IT – Key Ideas

## Task Generalization

**A wide variety of representative task to promote the task generalization**

| Capability | Domain | Task | IT Dataset | Size | Reference |
|---|---|---|---|---|---|
| Tasks | Finance | Relation Cls. | FingptFinred | 27,600 | Sharma et al. (2022) |
| | | NER | FingptNERCls | 13,500 | Yang et al. (2023) |
| | | | FingptNER | 511 | Alvarado et al. (2015) |
| | | Headline Cls. | FingptHeadline | 82,200 | Sinha et al. (2020) |
| | | Sentiment Cls. | SentimentCls | 47,600 | Yang et al. (2023) |
| | | | SentimentTra | 76,800 | Yang et al. (2023) |
| | | Summariz. | TradeTheEvent | 258,000 | Zhou et al. (2021) |
| IF/Chat | General | IF/Chat | SelfInstruct | 82,000 | Wang et al. (2022) |
| | | | SlimOrca | 518,000 | Lian et al. (2023) |
| | | | UltraChat | 774,000 | Ding et al. (2023) |
| | | | ShareGPT | 100,000 | Link |
| | Finance | QA | FinanceInstruct | 178,000 | Link |
| | | | FingptConvfinqa | 8,890 | Chen et al. (2022) |
| | | | FlareFinqa | 6,250 | Chen et al. (2021) |
| | | | FlareFiqa | 17,100 | Yang et al. (2023) |
| Reasoning | Math | QA | OrcaMath | 200,000 | Mitra et al. (2024) |
| | | | MetaMathQA | 395000 | Yu et al. (2023) |
| | | | MathInstruct | 262,000 | Yue et al. (2023) |
| | Code | QA | MagicodeInstruct | 111,000 | Luo et al. (2023) |
| | Finance | CFA Exam | Exercise | 2,950 | - |
| Total | | | | 3,161,401 | |

Ke, Ming, Joty - Adaptation of LLMs Tutorial, NAACL 2025

# IT – Key Ideas

## Training Workflow

Base Model ——→ **CPT+ IT** ————→ Next Stage

E.g., FinDAP

Base Model ——→ **CPT** ——→ **IT** ————→ Next Stage

E.g., FinLLM, FinTral (and many others)

FinDAP: Demystifying Domain-adaptive Post-training for Financial LLMs, Ke et al., 2025
FinTral: A Family of GPT-4 Level Multimodal Financial Large Language Models, Bhatia et al., 2024
FinLLM: Open-FinLLMs: Open Multimodal Large Language Models for Financial Applications, Huang et al., 2024

Ke, Ming, Joty - Adaptation of LLMs Tutorial, NAACL 2025

# IT – Key Ideas Summary

## Training Recipe

**Data Recipe:**
   **Synthetic data** (e.g., self-instruct)

**Model Recipe:**
   **Packing and Loss Mask**
   **Training Workflow** (e.g., CPT → IT, CPT+IT)

Synthetic data = text generated by LLM

## Seed Data

**Data Mixture:** A wide variety of representative to promote task generalization

**Data Budget ~** 1 Million

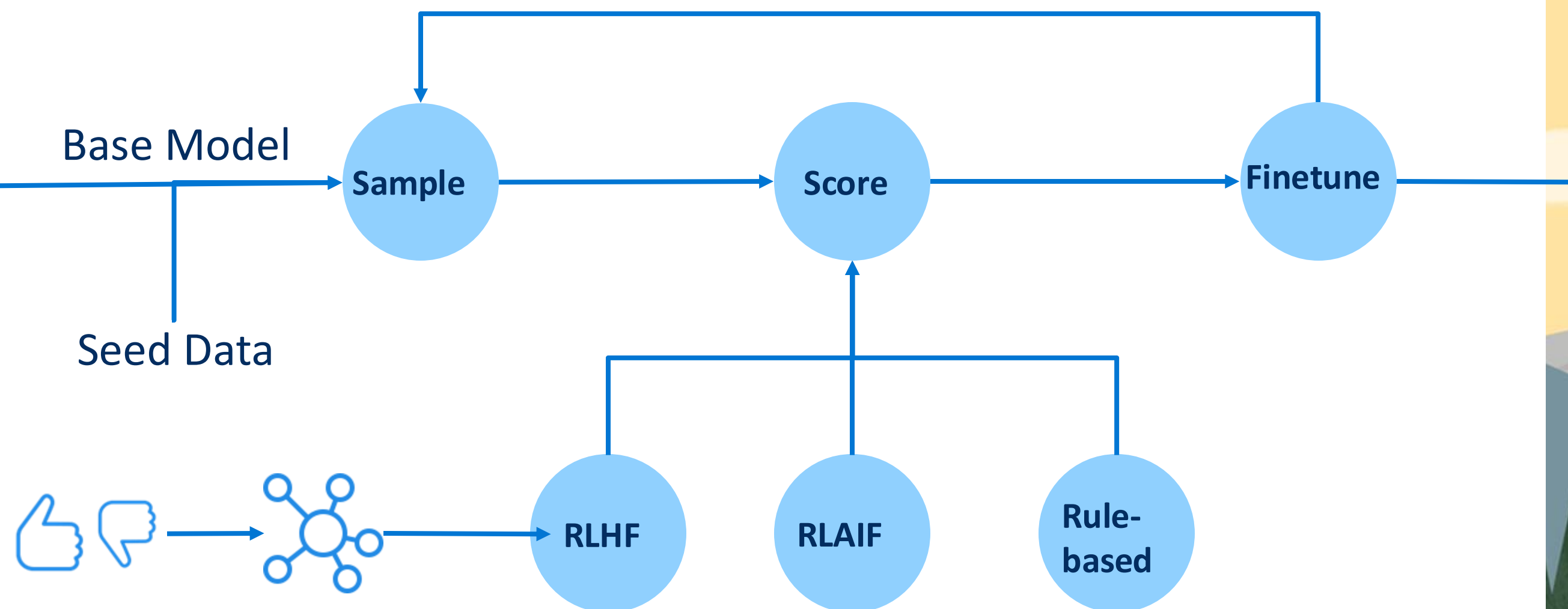# Supervised Preference Learning

# SPL – Role

## Style and Chat

Stronger training influence for style and chat capability

## More Capabilities

Continue building capabilities from instruction-tuned model, e.g., reasoning

Ke, Ming, Joty - Adaptation of LLMs Tutorial, NAACL 2025

# SPL – Example Workflow



Preference Learning Loop

Base Model

Sample → Score → Finetune

Seed Data

RLHF    RLAIF    Rule-based

Ke, Ming, Joty - Adaptation of LLMs Tutorial, NAACL 2025

# SPL – Key Considerations

## Training Recipe

**Data Recipe:** e.g., How to construct preference

**Model Recipe:**

**Algorithm**: How to optimize the preference reward?

**Training Workflow**: how to connect with other methods

## Seed Data

**Data Source:** Where to get the data?

**Data Mixture:** What should be included in the PL data?

**Data Budget:** How many data we need?

Ke, Ming, Joty - Adaptation of LLMs Tutorial, NAACL 2025

## DPO – Goal

$$\max_{\pi_\theta} \mathbb{E}_{x\sim\mathcal{D},y\sim\pi_\theta(y|x)}\left[r_\phi(x,y)\right] - \beta\mathbb{D}_{\mathrm{KL}}\left[\pi_\theta(y\mid x) \| \pi_{\mathrm{ref}}(y\mid x)\right]$$

**Optimize "reward" inspired
by human preferences**

**Constraint the model to not trust the
reward too much (preferences are
hard to model)**

**Main Questions:**

**1. How to implement the reward?**

**2. How to optimize the reward?**

# SPL – Key Ideas

## DPO – Preference / Reward modeling

**Chosen Completion**

**Prompt**

**Scores from optimal reward model**

$$p^*(y_1 \succ y_2 \mid x) = \frac{\exp\left(r^*(x, y_1)\right)}{\exp\left(r^*(x, y_1)\right) + \exp\left(r^*(x, y_2)\right)}.$$

**Rejected Completion**

**Key Idea:** **Probability $\propto$ Reward**

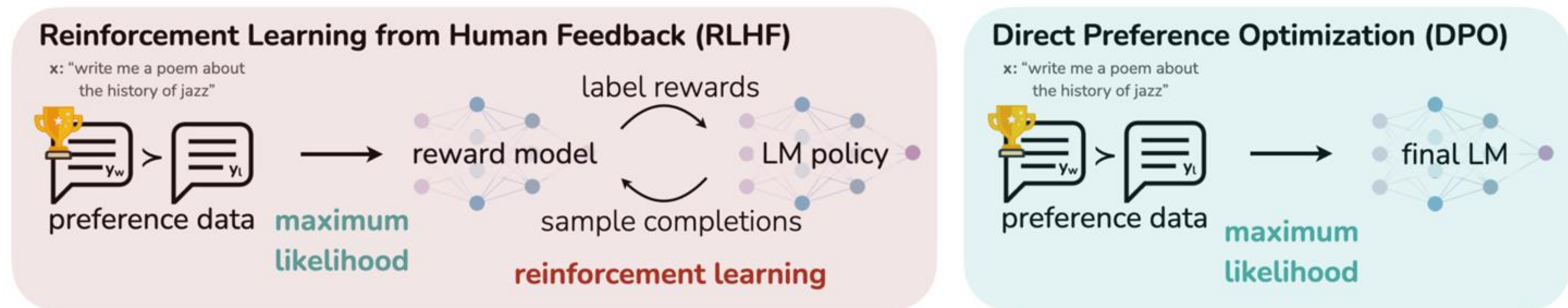**Obtaining point-wise Scalar reward of how good response is hard, but pairwise preference is easier and works!**

Ke, Ming, Joty - Adaptation of LLMs Tutorial, NAACL 2025

39

# SPL – Key Ideas

DPO

**If we just use gradient ascent on the equation**

With some math, we get: Direct Preference Optimization (DPO)



Direct Preference Optimization: Your Language Model is Secretly a Reward Model, Rafailov et al., 2023

# SPL – Key Ideas

## RLAIF

**Human Preferences (RLHF) vs. LLM-as-a-judge (RLAIF)**

Both source of preference data are used extensively

**In Frontier Labs:**

Human data used extensively as foundation

Synthetic data used to enhance behaviors (e.g., Constitutional AI)

**In Open Research:**

Synthetic data dominates (due to price)

Constitutional AI: Harmlessness from AI Feedbackl, Bai et al., 2022

# SPL – Key Ideas

## A Leading Synthetic Preference Method–UltraFeedback

**Key aspects**

Diverse model pool for completions

Diverse prompt pool

On-policy generations from checkpoints



UltraFeedback: Boosting Language Models with Scaled AI Feedback, Cui et al., 2024

# SPL – Key Ideas

Representative work with DPO – Zephyr, TuLU 70B….

**First model makes a splash with DPO**

**Fine-tune from Mistral 7b with UltraFeedback Datasets**

**Low learning rate (~5E-7) is good for DPO**



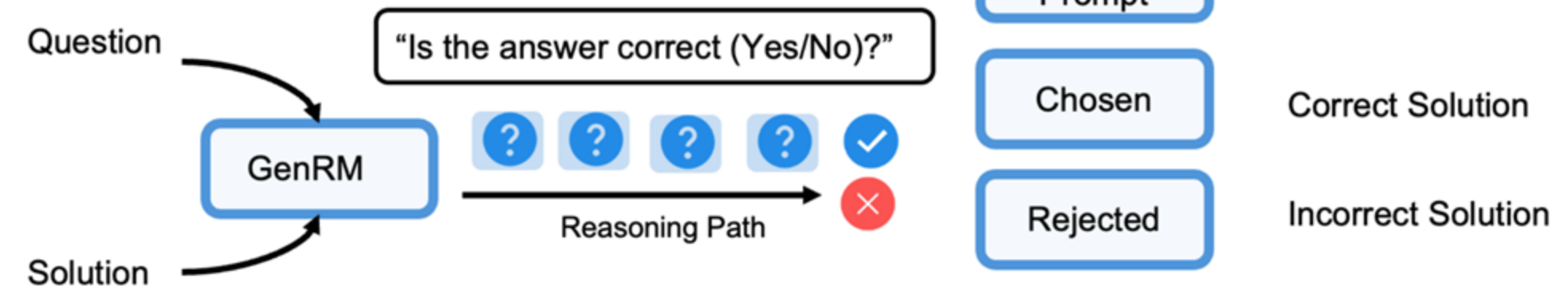Zephyr: Direct Distillation of LM Alignment, Tunstall, et al., 2023

# SPL – Key Ideas

## Synthesize Preference Data Focused on **Intermediate Preference**

**Final outcome preference**



Demystifying Domain-adaptive Post-training for Financial LLMs, Ke et al., 2025

# SPL – Key Ideas

## Synthesize Preference Data Focused on **Intermediate Preference**

**Final outcome preference**

**Intermediate outcome preference**

Identify and rectify the first erroneo step



Demystifying Domain-adaptive Post-training for Financial LLMs, Ke et al., 2025

# SPL – Key Ideas Summary

## Training Recipe

**Data Recipe:** Preference construction is often from diverse source (e.g., instruction pool, model pool) and cover fine-grained information (e.g., intermediate preference)

**Model Recipe:**

**Algorithm**: most popular: DPO

**Training Workflow**: usually after CPT and IT

## Seed Data

**Data Source:** often partial overlapping with IT

**Data Mixture:** Can be large scale (e.g., Math, Logic, Code, Science, Reasoning..)

**Data Budget:** ~ 1 million

Ke, Ming, Joty - Adaptation of LLMs Tutorial, NAACL 2025

# Coffee Break
# (30 Min)

# Reinforcement Learning

# RL – Role

## Beyond Human/AI Preference

RL as a training objective, learning from experience of interacting of the environment
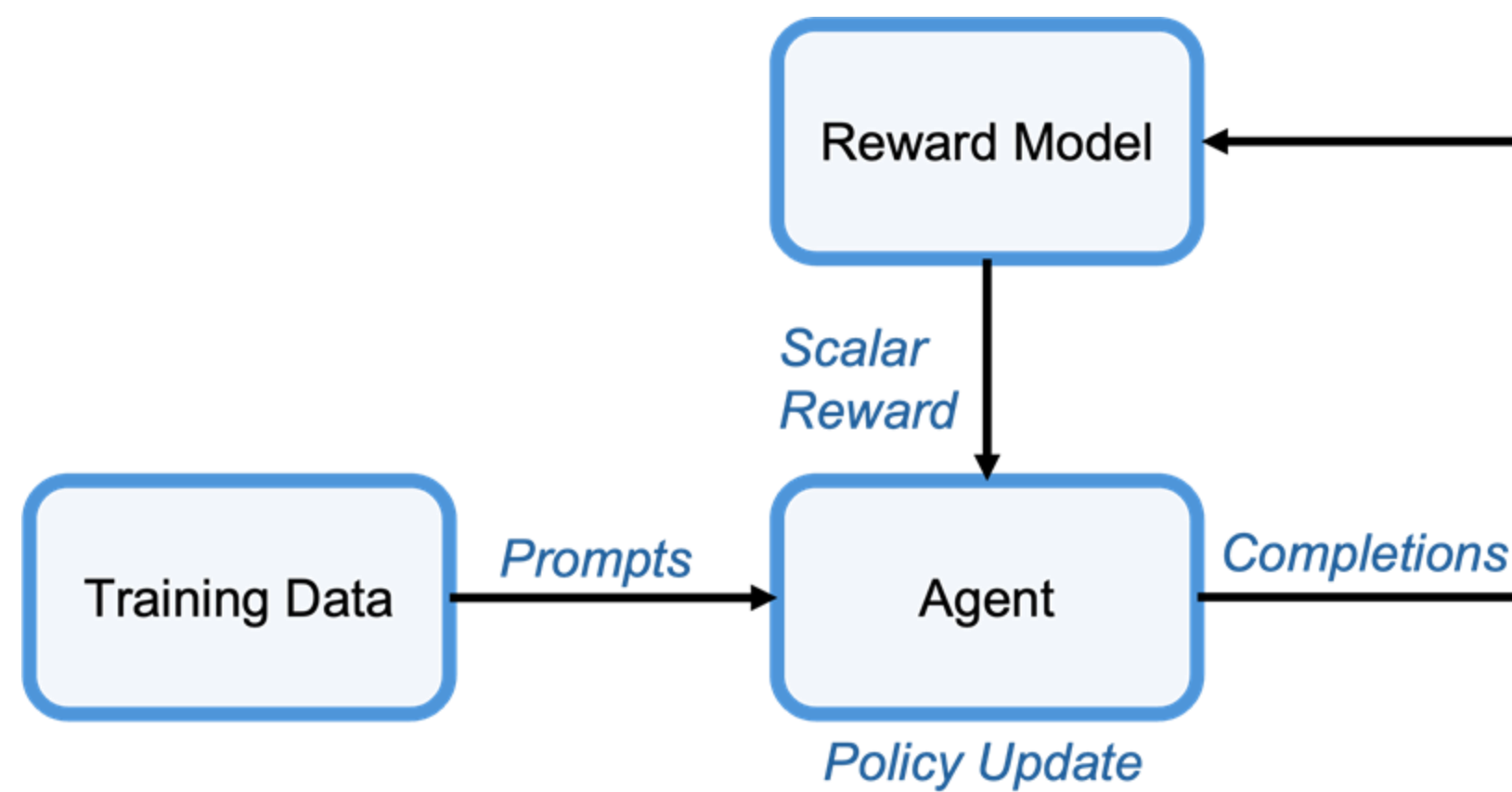
Recently show high-effectiveness

## Learn from Mistakes

RL methods naturally see both correct and a wide range of incorrect solutions.

This means they can:

improve targeted capabilities **without** degradation on other out-of-domain capabilities

Ke, Ming, Joty - Adaptation of LLMs Tutorial, NAACL 2025

# RL – Example Workflow



Ke, Ming, Joty - Adaptation of LLMs Tutorial, NAACL 2025

# RL – Key Considerations

## Training Recipe

**Model Recipe:**

**Algorithm**: How to optimize the reward effectively and efficiently?

**Training Workflow**: how to connect with other methods

## Seed Data

**Data Source:** Where to get the data?

**Data Mixture:** What should be included in the RL data?

**Data Budget:** How many data we need?

Ke, Ming, Joty - Adaptation of LLMs Tutorial, NAACL 2025

# RL – Key Ideas

## From DPO to RL

$$\max_{\pi_\theta} \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_\theta(y|x)} \left[ r_\phi(x, y) \right] - \beta \mathbb{D}_{\text{KL}} \left[ \pi_\theta(y \mid x) \mid\mid \pi_{\text{ref}}(y \mid x) \right]$$

**Optimize "reward" inspired by human preferences**

**Constraint the model to not trust the reward too much (preferences are hard to model)**
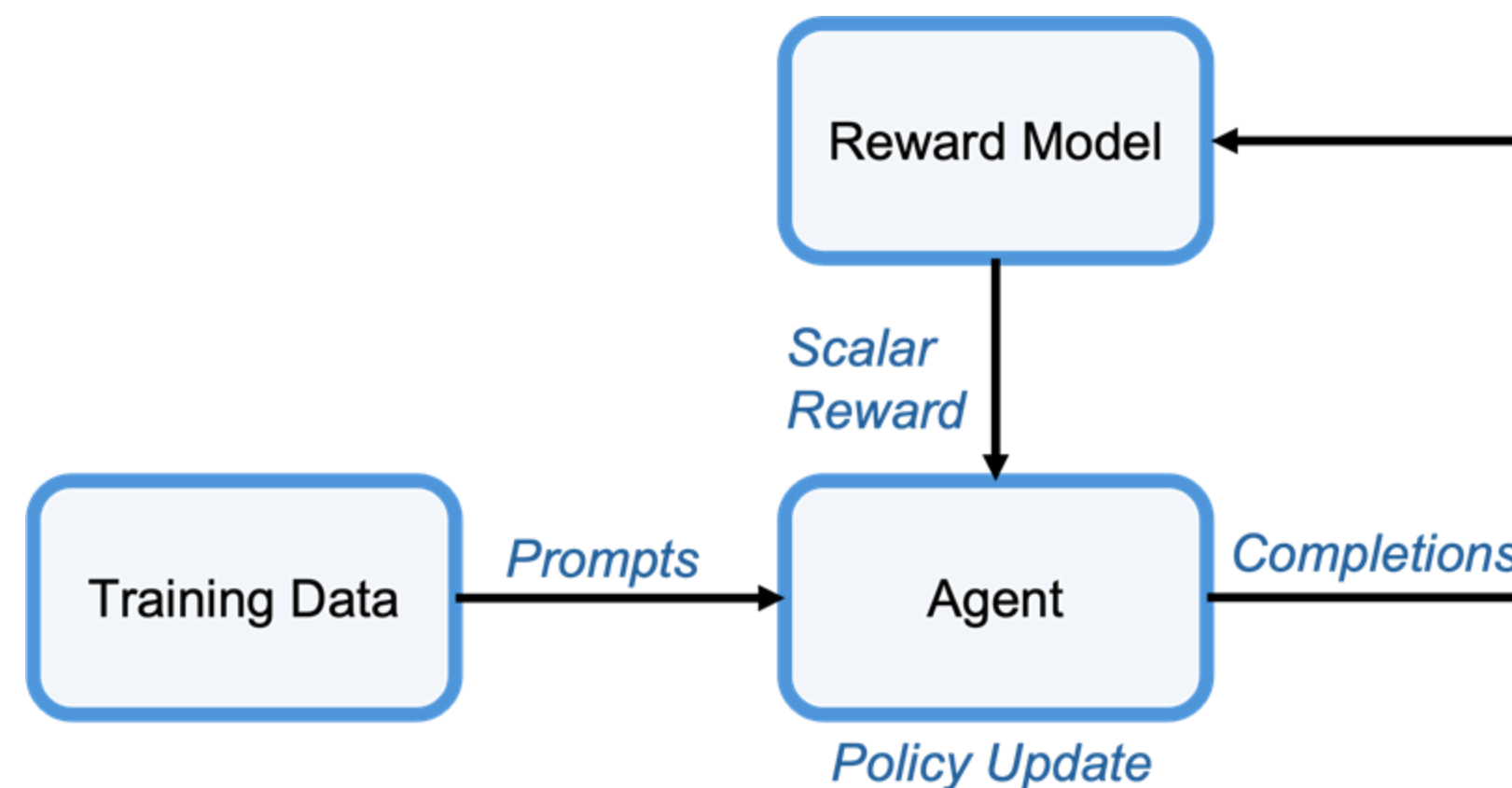
**Main Questions:**

1. **How to implement the reward?**

2. **How to optimize the reward?**

Ke, Ming, Joty - Adaptation of LLMs Tutorial, NAACL 2025

# RL – Key Ideas

## From DPO to RL

**What if we choose not to use pairwise preference but still rely on scalar reward**
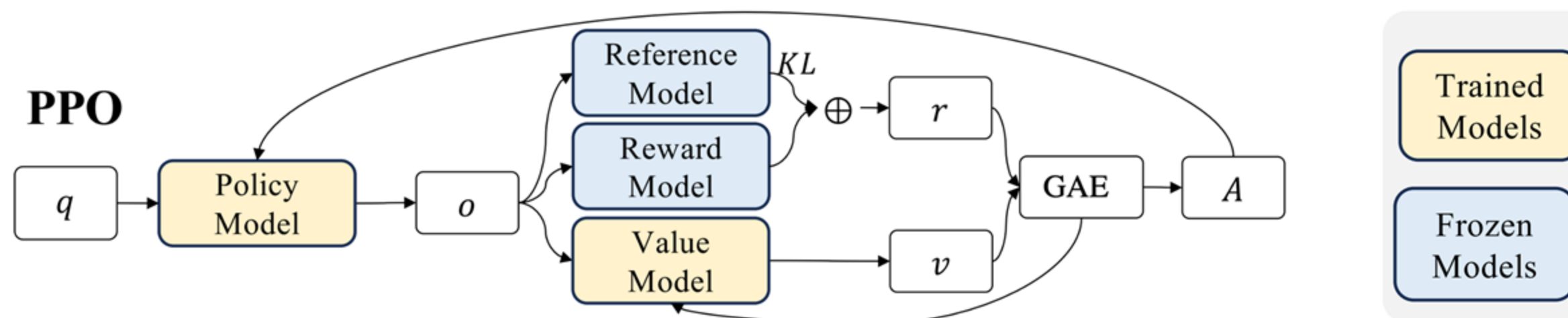


Ke, Ming, Joty - Adaptation of LLMs Tutorial, NAACL 2025

PPO

**One popular method is PPO**

**(effective but expensive: 4 copies of model)**



Proximal Policy Optimization Algorithms

John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, Oleg Klimov
OpenAI
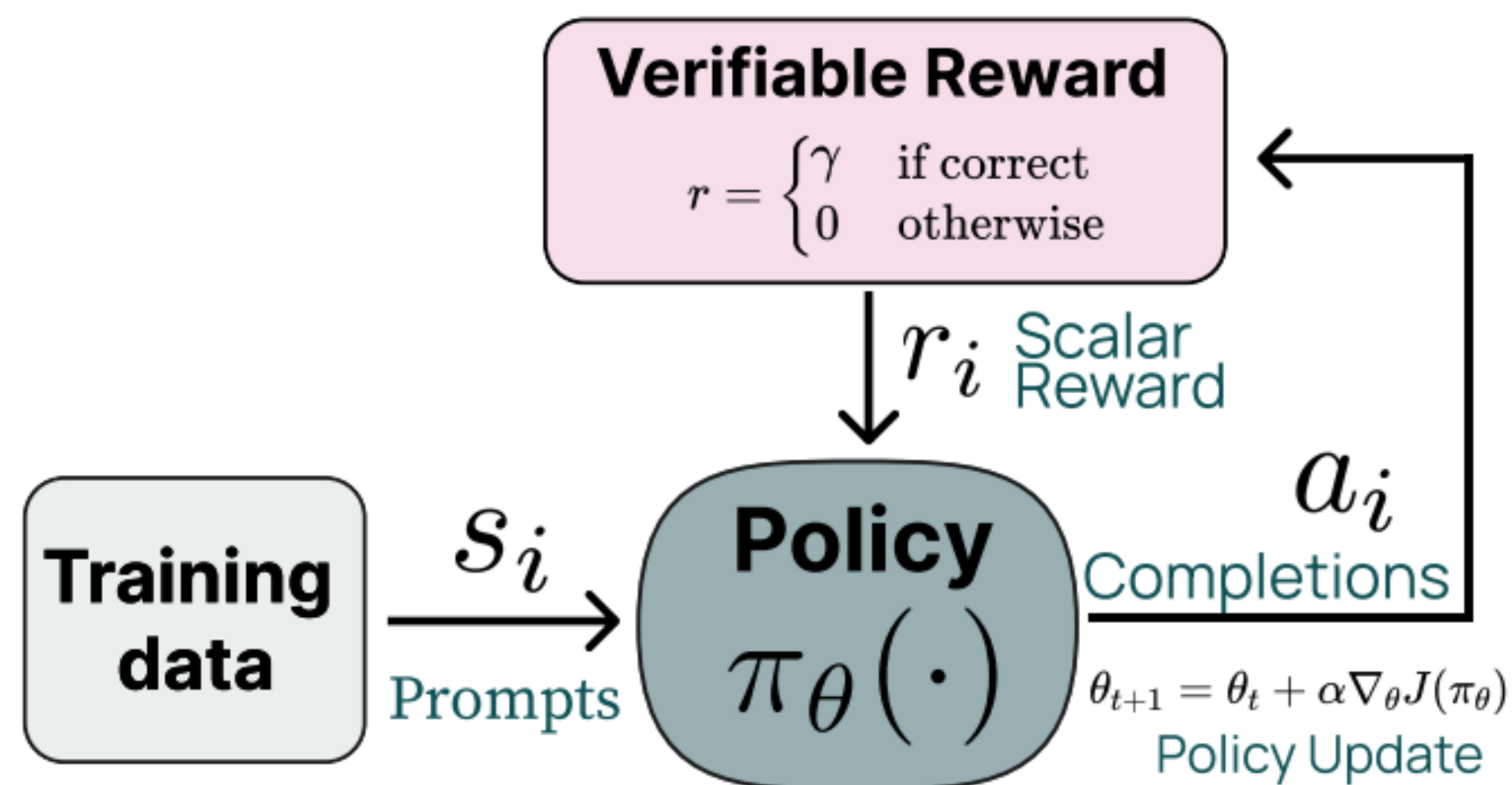{joschu, filip, prafulla, alec, oleg}@openai.com

Ke, Ming, Joty - Adaptation of LLMs Tutorial, NAACL 2025

# RL – Key Ideas
## RL with Verifiable Reward (RLVR)

**Since the scalar reward is hard to get, one method is to use verifiable reward (e.g., math)**

Reward model is also eliminated



**Verifiable Reward**

$$r = \begin{cases} \gamma & \text{if correct} \\ 0 & \text{otherwise} \end{cases}$$

$r_i$ Scalar Reward

**Training data** $\xrightarrow{s_i}$ **Policy** $\pi_\theta(\cdot)$

Prompts

$a_i$ Completions

$\theta_{t+1} = \theta_t + \alpha \nabla_\theta J(\pi_\theta)$

Policy Update

Tülu 3: Pushing Frontiers in Open Language Model Post-Training, Lambert et al., 2025

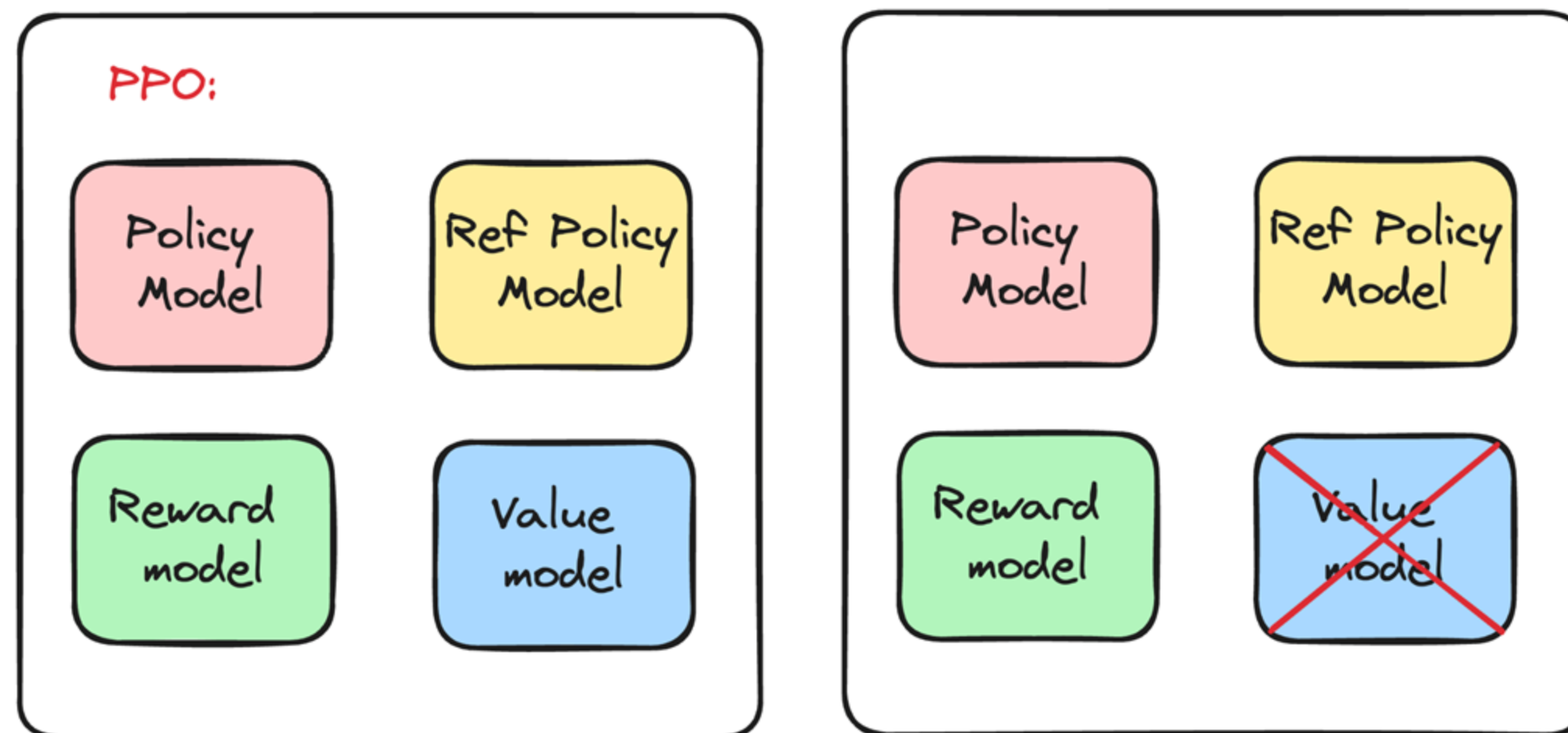Ke, Ming, Joty - Adaptation of LLMs Tutorial, NAACL 2025

# RL – Key Ideas

## Can We Get Rid of the Value Model?

**But this is still limited, can we get rid of the value model?**

The answer to this question leads to many RL algorithm variants for LLM

Ke, Ming, Joty - Adaptation of LLMs Tutorial, NAACL 2025

# RL – Key Ideas

## Can We Get Rid of the Value Model?

**Core Trick**

**Value Model =** a model (LLM) that estimates the baseline expected return at each time step (token), so we can measure how much better or worse the actual outcome was compared to this expectation (this difference is called advantage).

# RL – Key Ideas

## Can We Get Rid of the Value Model?

**Core Trick**

***But,*** *do we need we really need to figure out which* **token** *made the reader happy?*

*Can we just ask "Is the answer good?" If yes → reinforce. No need to slice the blame*

> **Key Innovation:**
>
> **Value attributed to each token → group of tokens (e.g., full response)**
>
> **Now the value is directly tie to the reward, no value model required to estimate expected return at each time step.**
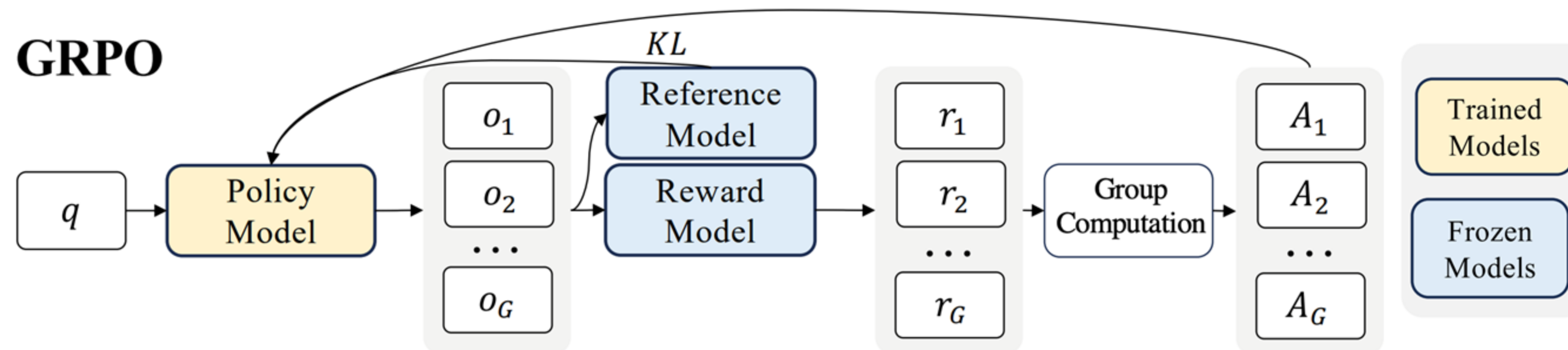
GRPO

**Action = full response**

**Advantage = Preference ranking across a group**



DeepSeekMath: Pushing the Limits of Mathematical Reasoning in Open Language Models

Ke, Ming, Joty - Adaptation of LLMs Tutorial, NAACL 2025

# RL – Key Ideas

Another RL Variant: RLOO

**Action = full response**

**Advantage = Leave-One-Out reward baseline**

$$A = R(x, y) - \frac{1}{n-1} \sum_{j \neq i} R(x, y_j)$$

**Reward for the current response**

**All other responses in the batch**

## Back to Basics: Revisiting REINFORCE Style Optimization for Learning from Human Feedback in LLMs

**Arash Ahmadian**
*Cohere For AI*

**Chris Cremer**
*Cohere*

**Matthias Gallé**
*Cohere*

Ke, Ming, Joty - Adaptation of LLMs Tutorial, NAACL 2025

60

# RL – Key Ideas Summary

## Training Recipe

**Model Recipe:**
    **Algorithm**: Value model is eliminated by taking group of token as action and define advantage based on those group of tokens (various across RL algorithms. It is still an active research topic)

    **Training Workflow**: usually serve as the last method in the workflow (e.g., after CPT, IT, and PL)

## Seed Data

**Data Source:** often partial overlapping with IT

**Data Mixture:** Can be large scale (e.g., Math, Logic, Code, Science, Reasoning..)

**Data Budget ~** 10 thousand (recent research shows that even a small amount, even just 1-shot can make a different. Still actively research)

Ke, Ming, Joty - Adaptation of LLMs Tutorial, NAACL 2025