



# Domain-adaptive Post-training For Financial LLMs

Zixuan Ke

<https://vincent950129.github.io/>



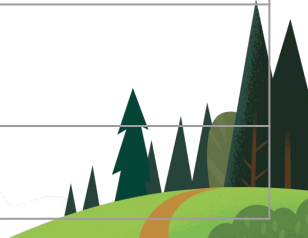
# FinDAP

## Highlighted outcomes



Resulting model  
(Llama-Fin-8b), a small but  
mighty Finance LLM!

Benchmark	Llama-Fin (Ours, 8b)	GPT4o
FPB (Financial Sentiment analysis)	<u>91.13</u>	82.16
FiQA SA (Financial Sentiment analysis)	<u>95.32</u>	68.51
NER (Financial Named Entity Recognition)	<u>76.69</u>	43.02
EDTSUM (Financial Abstractive Summarization)	<u>53.78</u>	18.15
Finance Bench (Financial Open QA)	<u>54.00</u>	51.30
SM-Bigdata (Stock Movement Prediction)	<u>54.14</u>	49.18
Flare-German (Credit Scoring)	<u>64.00</u>	17.00
.....		



# FinDAP

## Outcome



Evaluation set is open source

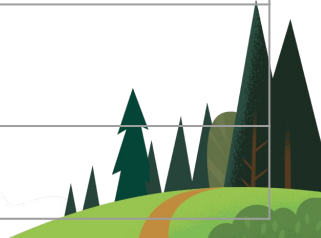
Dataset Viewer: Salesforce/FinEval

Subset (25) CFChallenge · 90 rows

query string	answer string	source string
Carefully read the scenario provided and the subsequent question. Your task is to analyze the scenario and select the most appropriate answer from the options A, B and C. Scenario: TOPIC...	A	sample_test
Carefully read the scenario provided and the subsequent question. Your task is to analyze the scenario and select the most appropriate answer from the options A, B and C. Scenario: TOPIC...	A	sample_test
Carefully read the scenario provided and the subsequent question. Your task is to analyze the scenario and select the most appropriate answer from the options A, B and C. Scenario: TOPIC...	A	sample_test
Carefully read the scenario provided and the subsequent question. Your task is to analyze the scenario and select the most appropriate answer from the options A, B and C. Scenario: TOPIC...	C	sample_test
Carefully read the scenario provided and the subsequent question. Your task is to analyze the scenario and select the most appropriate answer from the options A, B and C. Scenario: TOPIC...	A	sample_test
Carefully read the scenario provided and the subsequent question. Your task is to analyze the scenario and select the most appropriate answer from the options A, B and C. Scenario: TOPIC...	B	sample_test
Carefully read the scenario provided and the subsequent question. Your task is to analyze the scenario and select the most appropriate answer from the options A, B and C. Scenario: TOPIC...	C	sample_test

<https://huggingface.co/datasets/Salesforce/FinEval>

Benchmark	Llama-Fin (Ours, 8b)	GPT4o
FPB (Financial Sentiment analysis)	<u>91.13</u>	82.16
FiQA SA (Financial Sentiment analysis)	<u>95.32</u>	68.51
NER (Financial Named Entity Recognition)	<u>76.69</u>	43.02
EDTSUM (Financial Abstractive Summarization)	<u>53.78</u>	18.15
Finance Bench (Financial Open QA)	<u>54.00</u>	51.30
SM-Bigdata (Stock Movement Prediction)	<u>54.14</u>	49.18
Flare-German (Credit Scoring)	<u>64.00</u>	17.00
.....		



# Research Question



Given a strong general-purpose LLM (*e.g.*, Llama3-8b-inst), how to effectively adapt it to a target domain (*e.g.*, finance) by post-training? What criteria are desirable for successful adaptation? What are effective training recipes with respect to data and model?



# Research Question



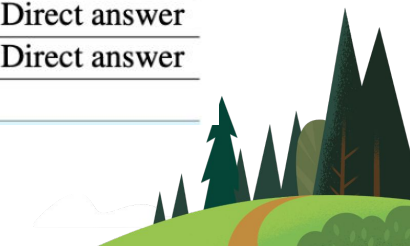
Answer from related work:  
(e.g., PIXIU, FinLLM, FinTral, Palmyra-Fin, FinMa, Finance-LLM, FinLLaVA)\*

*Follow standard methods:*

Continual Pre-training (CPT) →  
Instruction-tuning (IT) → Preference  
Alignment (PA)

Given a strong general-purpose LLM (e.g., Llama3-8b-inst), how to effectively adapt it to a target domain (e.g., finance) by post-training? What criteria are desirable for successful adaptation? What are effective training recipes with respect to data and model?

Finance LLM	Capabilities	Recipe		Evaluation
		Model Recipe	Data Recipe	
AdaptLLM	Concept	CPT	<b>CPT:</b> Financial text + heuristic QAs constructed from the text	Financial + Classification tasks + Direct answer
PIXIU	Task	IT	<b>IT:</b> Financial tasks	Financial + Classification tasks + Direct answer
FinLLM	Concept, Task	CPT → IT	<b>CPT:</b> Financial text + Fineweb; <b>IT:</b> Filtered Financial tasks	Financial + Classification tasks + Direct answer
FinTral	Concept, Task	CPT → IT → PA	<b>CPT:</b> Financial text; <b>IT:</b> Financial tasks; <b>PA:</b> Outcome signal only	Financial + Classification tasks + Direct answer
Palmyra-Fin			SoTA public checkpoint, but recipe is not disclosed	



# Research Question



Answer from related work:  
(e.g., PIXIU, FinLLM, FinTral, Palmyra-Fin  
etc.)\*

*Follow standard methods:*

Continual Pre-training (CPT) →  
Instruction-tuning (IT) → Preference  
Alignment (PA)

## **This work (FinDAP)**

This is not enough! Domain-adaptive  
post-training is unique to pre-training and  
general post-training and we need a systematic  
and principle approach

Given a strong general-purpose LLM (e.g., Llama3-8b-inst), how to effectively adapt it to a target domain (e.g., finance) by post-training? What criteria are desirable for successful adaptation? What are effective training recipes with respect to data and model?



# FinDAP



## Factors to consider

- **Desirable capabilities** for the target domain (e.g., reasoning...)
- **Training Recipe**
  - Original pre-trained LLM already possess strong general capabilities and knowledge
    - Catastrophic Forgetting
    - Knowledge Transfer
  - Construct preference data for reasoning in preference alignment (PA)
- **Implementation of the recipe** (training data)
  - Quantity vs. Quality
    - Literature found that small amount of general data is enough to mitigate forgetting
    - While learning domain-specific knowledge typical require more data
- **Evaluation**
  - Different capabilities may require different evaluation methods
    - e.g., reasoning tasks may want CoT evaluation

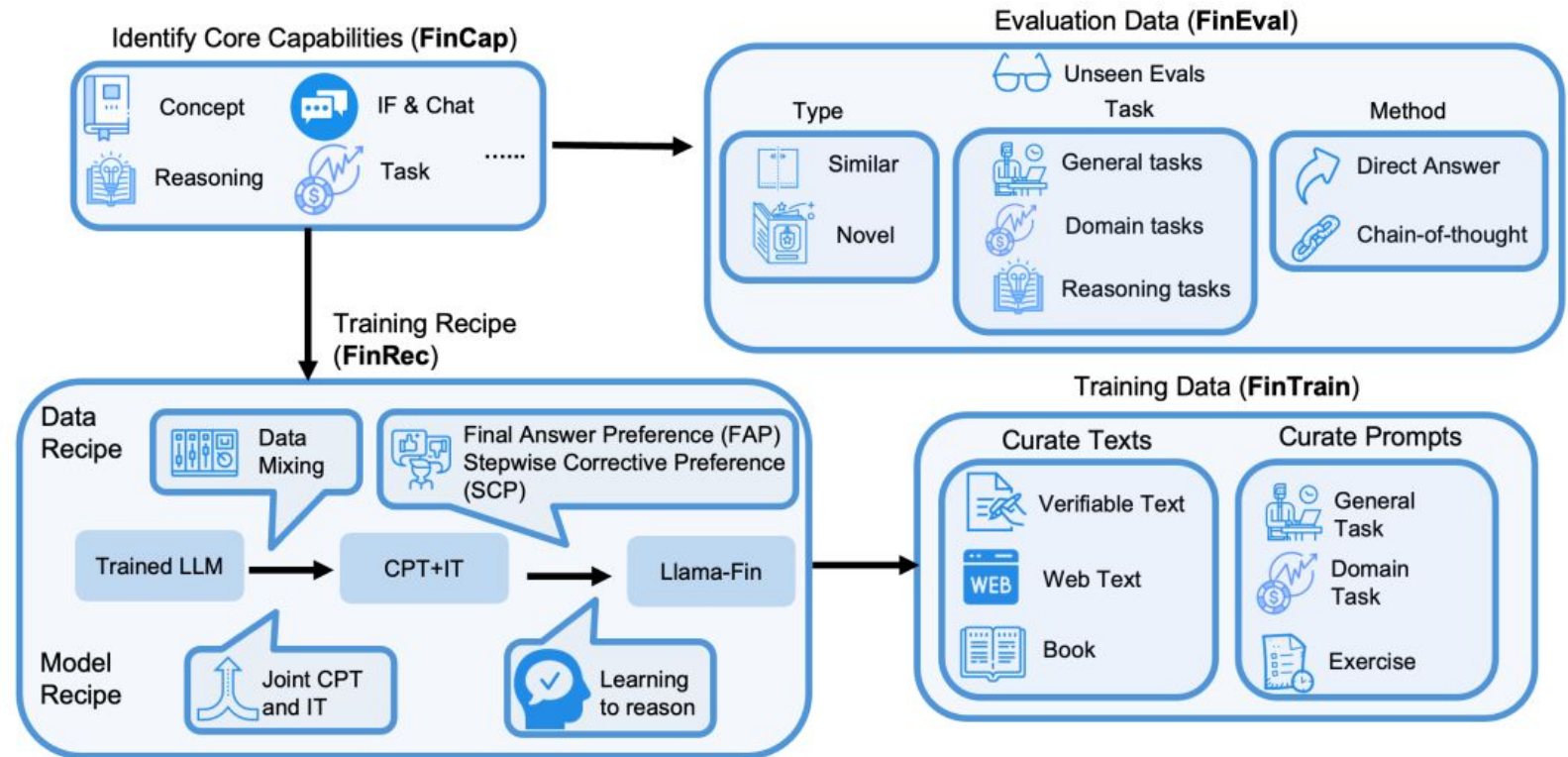


# FinDAP

## Our Framework



- **FinCap:** Core capacities required for finance domain
- **FinRec:** Our training Recipe
- **FinTrain:** a curated set of training datasets implement FinRec
- **FinEval:** A comprehensive evaluation framework



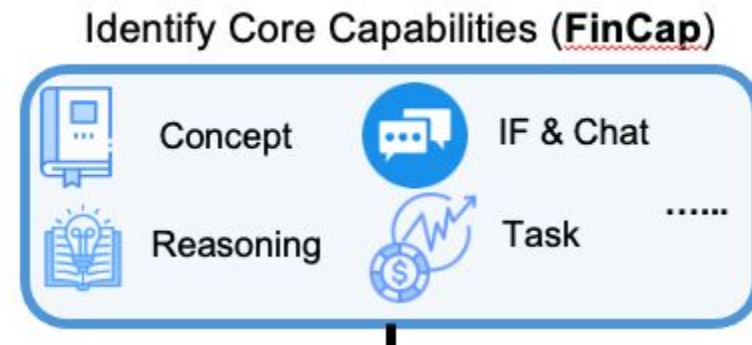


# FinDAP



## FinCap: Core capacities required for finance domain

- We propose 4 main capabilities based on the fundamental requirement in FinAI
  - Understanding **domain-specific concepts** to process financial language accurately, performing **domain-specific tasks** to solve real-world problems, **reasoning** effectively to analyze complex financial data, and **following instructions** to interact naturally in practical applications.
  -
- **Domain-specific concepts**
- **Domain-specific tasks**
- **Instruction-following**
- **Reasoning**

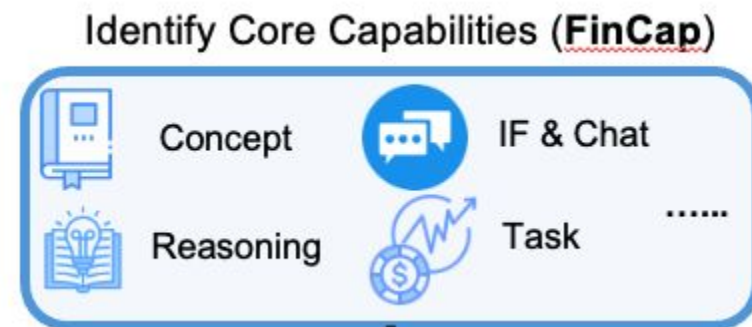


# FinDAP



## FinCap: Core capacities required for finance domain

- We propose 4 main capabilities based on the fundamental requirement in FinAI
  - Understanding **domain-specific concepts** to process financial language accurately, performing **domain-specific tasks** to solve real-world problems, **reasoning** effectively to analyze complex financial data, and **following instructions** to interact naturally in practical applications.



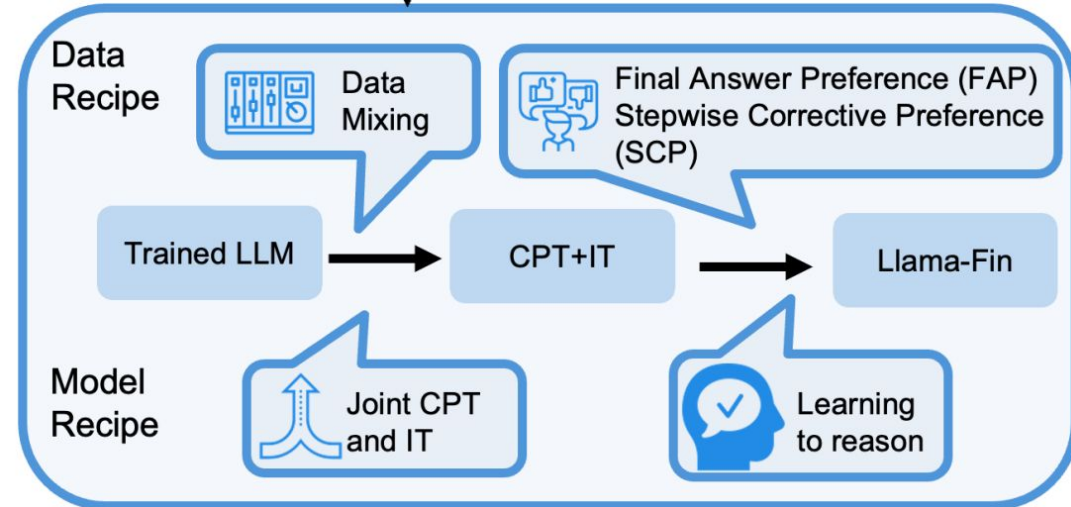
- 
- Domain-specific concepts → Continual Pre-training (CPT)
- Domain-specific tasks → Instruction Tuning (IT)
- Instruction-following → Instruction Tuning (IT)
- Reasoning → Preference Alignment (PA)



# FinDAP

## FinRec: Our training Recipe

- Model Recipe
  - Joint training CPT and IT
    - Why?
      - CPT alone causes forgetting on instruction-following abilities.
      - A joint training can further improve generalization
        - Concepts are often inherently more generalizable due to the shared nature of concepts across tasks
      - Implementation
        - Since the only different is whether to mask-out the instruction, we can simply mixing up their data to achieve jointly training
        - CPT data size is usually larger, we down-sample it to match the IT size
    - PA for reasoning tasks
      - Assign higher probability mass to better generations, has been shown to be effective in enhancing reasoning capabilities of LLMs
      - Employ DPO (detailed in data recipe)



# FinDAP

## FinRec: Our training Recipe

- Data Recipe

- In-domain, general-domain and mixture

- Most FinLLMs use in-domain data only

- This exclusive reliance on in-domain data can lead to forgetting of general knowledge in the original pre-trained LLM.

- We conduct systematic investigation ( $\{CPT/IT/PA\}-\{In/Gen/Mix\}$ )

- CPT

- While CPT-In and CPT-Gen outperforms in financial and general tasks, respectively, CPT-Mix achieves the best → mixing data sources effectively mitigates forgetting of general knowledge

- IT

- IT-Mix slightly outperforms than other data versions → mixing general tasks remains helpful to mitigating forgetting of general concepts and tasks, although the effect is much less pronounced compared to CPT.

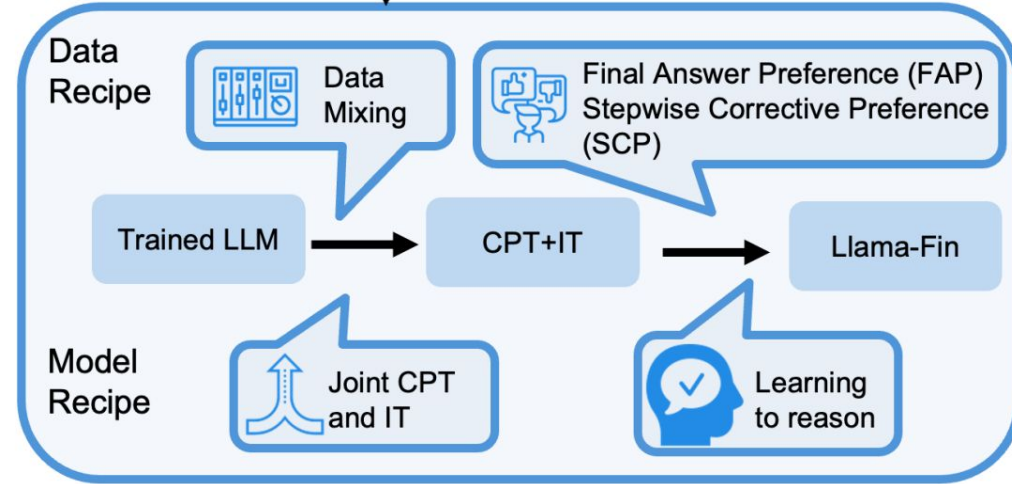
- PA.

- PA-In performs comparably to PA-Mix, indicating that it is NOT essential to include general tasks to prevent forgetting of concepts or tasks, unlike the cases of CPT and IT.

Forgetting  
impact  
decrease



Mixture of in- and general-domain data for CPT+IT



# FinDAP

## FinRec: Our training Recipe

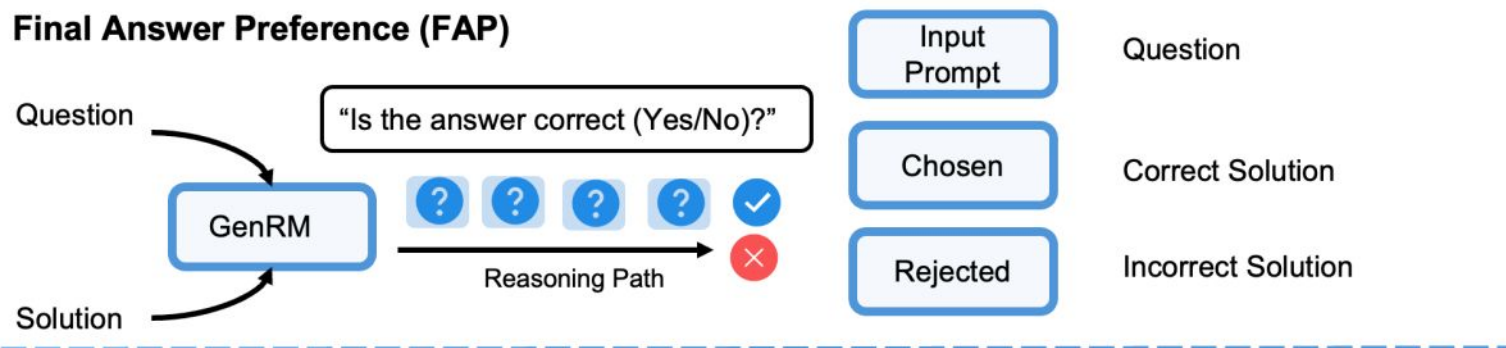
- Data Recipe
  - Preference data with outcome and process signals/reward
    - Two Regimes: Learning to reason (DS-R1...) vs inference scaling (OAI-O1..)
      - We adopt “learning to reason” as finance domain often require rapid responses
    - Learning to reason
      - Trajectories collection
        - Search-based
          - RM/verifier to guide the search
        - Revision-based
          - Iterative refinement
        - We adopt the **search-based method** as revision-based shows mixed results and have not yet been well established as reliable for achieving improvements
      - Training from trajectories (DPO in this work)



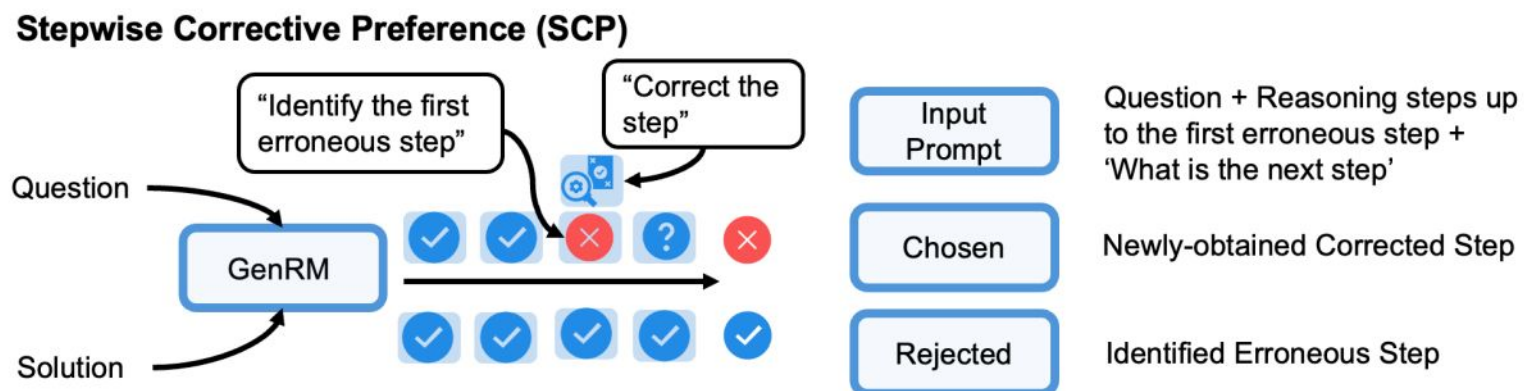
## FinRec: Our training Recipe

- Learning to reason: Search-based trajectories collection
  - Reward model / verifier
    - We employ generative RM with strong pre-trained LLM (i.e., GPT4o)

Final answer  
(Outcome) Level



Reasoning steps  
(Process) Level

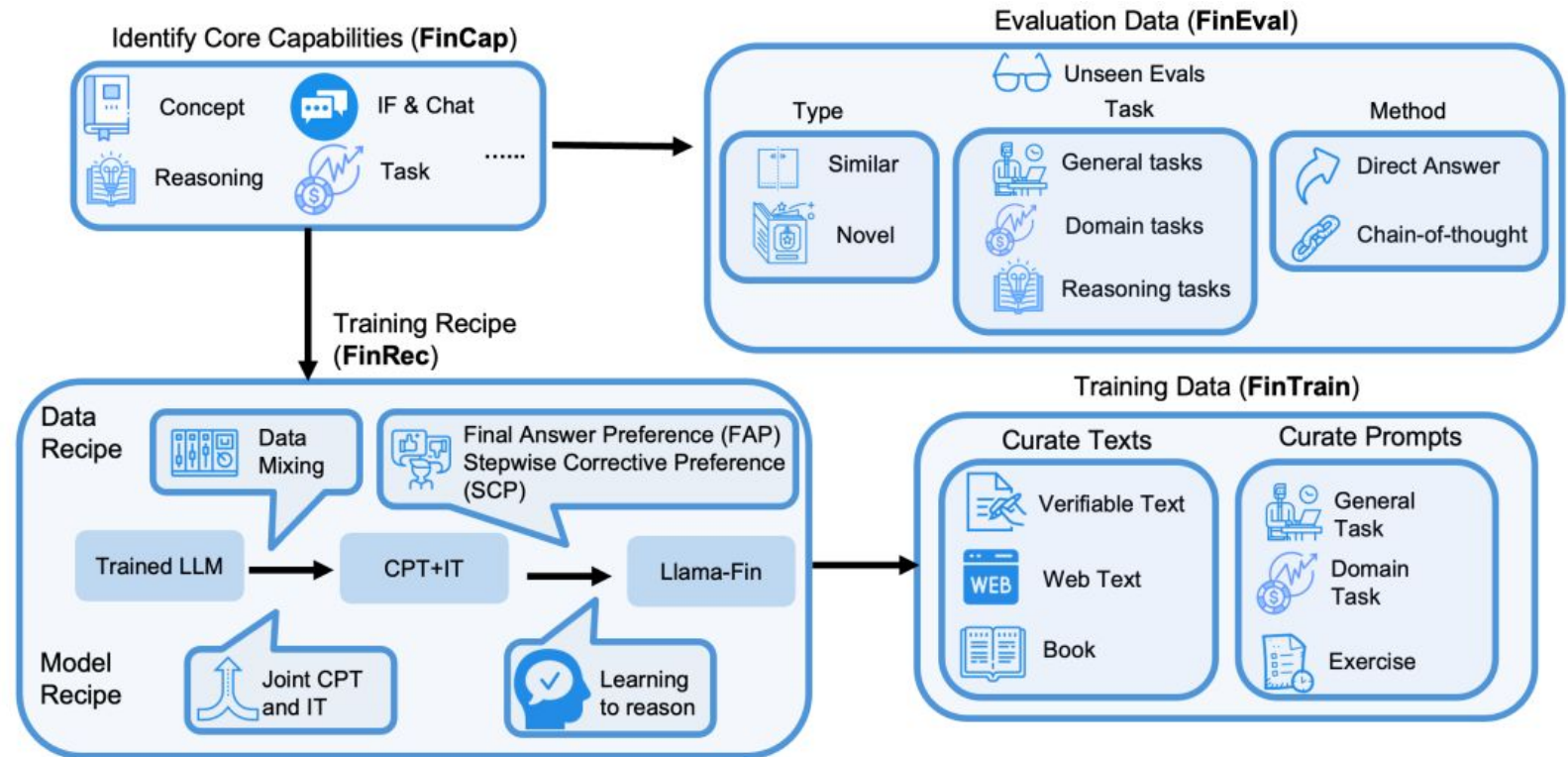


# FinDAP

## Our Framework



- **FinCap:** Core capacities required for finance domain
- **FinRec:** Our training Recipe
- **FinTrain:** a curated set of training datasets implement FinRec
- **FinEval:** A comprehensive evaluation framework

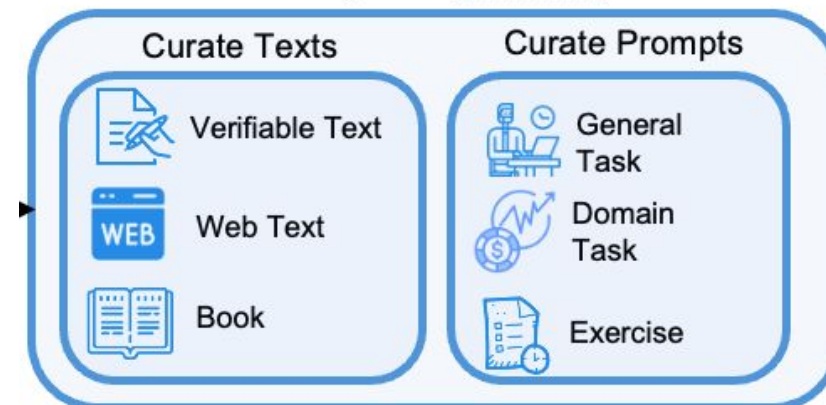


# FinDAP

FinTrain: a curated set of training datasets implement FinRec

- **General text (used in CPT)**
  - **Goal:** mitigate forgetting
  - **Literature:** a ‘small’ amount of general text (as little as 1%) can effectively mitigate the forgetting issue
  - **FinDAP:** focus on collecting a relatively small but high-quality set of general-domain text.
  - **Verifiable:** text written by humans and previously used in supervised task
- **Finance text (used in CPT)**
  - **Goal:** domain-specific knowledge
  - Diverse and large-scale:
    - Web (URLs based filtering)
    - Books

## Training Data (FinTrain)



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 (2019)
		DialogueStudio	1,070,000	Zhang et al. (2023)
Finance	Fineweb-Fin	Book-Fin	4,380,000	-
		Book-Fin	4,500	-
<b>Total</b>			10,177,294	

Table 3: Summary of curated texts. New datasets released with FINDAP are color-highlighted for emphasis.

CPT datasets totally ~ 6B tokens

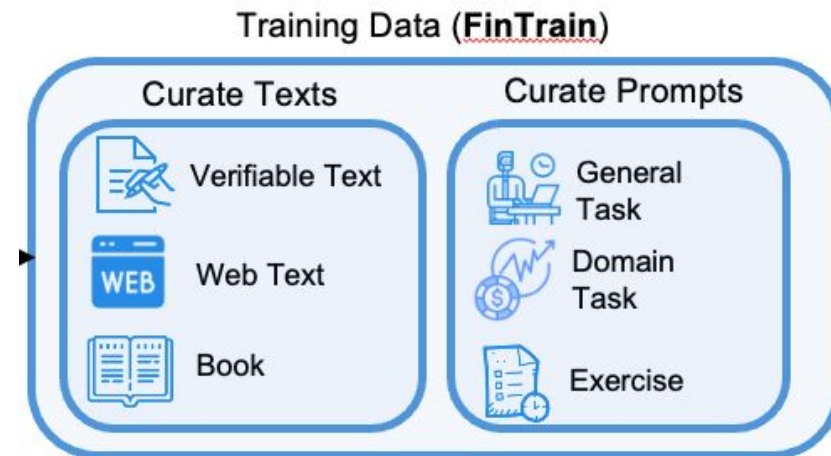




# FinDAP

FinTrain: a curated set of training datasets implement FinRec

- Prompts (used IT and PA)
  - Corresponding to each capabilities
  - Diversity
  - Previously shown perwell well (e.g., UltraQA)
  - Potential reasoning steps provided (e.g., Exercise)



Capability	Domain	Task	IT Dataset	Size	Reference		
Tasks	Finance	Relation Cls. NER	FingptFinred	27,600	Sharma et al. (2022)		
			FingptNERCls	13,500	Yang et al. (2023)		
			FingptNER	511	Alvarado et al. (2015)		
		Headline Cls. Sentiment Cls.	FingptHeadline	82,200	Sinha et al. (2020)		
			SentimentCls	47,600	Yang et al. (2023)		
			SentimentTra	76,800	Yang et al. (2023)		
IF/Chat	General	Summariz. IF/Chat	TradeTheEvent	258,000	Zhou et al. (2021)		
			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)
					OrcaMath	200,000	Mitra et al. (2024)
Reasoning	Math	QA	MetaMathQA	395,000	Yu et al. (2023)		
			MathInstruct	262,000	Xiang Yue (2023)		
			MagicodeInstruct	111,000	Luo et al. (2023)		
			Finance	CFA Exam	<b>Exercise</b>	2,950	-
<i>Total</i>				3,161,401			

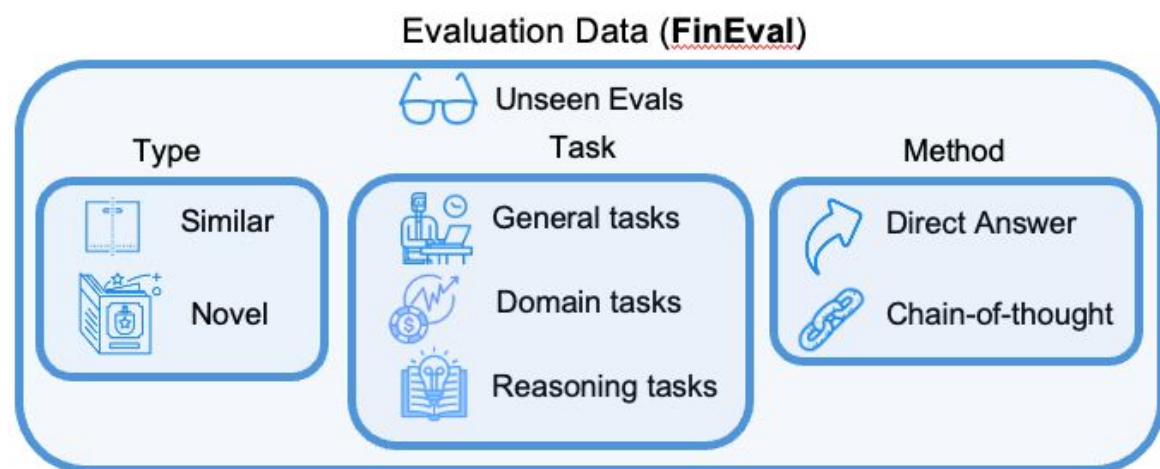
Table 4: Summary of our curated prompts. New datasets released with FINDAP are color-highlighted for emphasis. For datasets without formal references but only a URL, we provide their links.



# FinDAP

## FinEval: A comprehensive evaluation framework

- **Types**
  - Similar (to training)
    - The task type has been seen
    - **Goal:** outperform the best model (GPT4o..), even if we are small (post-training from LLama3-8b-instruct)
  - Novel
    - A new task type
    - **Goal:** outperform our original pre-train LLM (LLama3-8b-instruct)
- **Tasks**
  - Corresponding to the 4 capabilities
- **Methods**
  - CoT for reasoning tasks



Capability	Domain	Task	Evaluation Dataset	Size	Reference
<b>Unseen - Similar</b>					
Tasks	Finance	Sentiment Analysis	FPB	970	Malo et al. (2014)
			FiQA SA	235	Maia et al. (2018)
		Monetary policy Stance	FOMC	496	Shah et al. (2023)
		Named entity recognition	NER	98	Alvarado et al. (2015)
		Abstractive Summarization	EDTSUM	2,000	Zhou et al. (2021)
<i>Total</i>			3,799		
<b>Unseen - Novel</b>					
Concept	General	Knowledge Recall	MMLU	14,042	(Hendrycks et al., 2021)
			AI2-ARC	3,548	Clark et al. (2018)
			Nq-open	7,842	Kwiatkowski et al. (2019)
			<b>MMLU-Finance</b>	1,460	-
Tasks	Finance	Extractive Summarization	Flare-ECTSUM	495	Mukherjee et al. (2022)
			MLESG	300	Chen et al. (2023b)
		Rumour Detection	MA	500	Yang et al. (2020)
		Stock Movement Prediction	SM-Bigdata	1,470	Soun et al. (2022)
			SM-ACL	3,720	Xu and Cohen (2018)
			SM-CIKM	1,140	Wu et al. (2018)
		Fraud Detection	CRA-CCF	2,280	Feng et al. (2024)
			CRA-CCFraud	2,100	Feng et al. (2024)
		Credit Scoring	Flare-German	200	Hofmann (1994)
			Flare-Australian	139	Quinlan (1987)
			CRA-LendingClub	2,690	Feng et al. (2024)
		Distress Identification	CRA-Polish	1,740	Feng et al. (2024)
			CRA-Taiwan	1,370	Feng et al. (2024)
			CRA-ProroSeguro	2,380	Feng et al. (2024)
		Claim Analysis	CRA-TravelInsurance	2,530	Feng et al. (2024)
Flare-TATQA	1,670		Zhu et al. (2021)		
IF/Chat Reasoning	General	Open QA	Finance Bench	150	Islam et al. (2023)
		Precise IF	MT-bench	80	Zheng et al. (2023)
	Math	Reasoning	MathQA	2,985	Amini et al. (2019a)
		Social Reasoning	Social-IQA	2,636	Welbl et al. (2017)
	General	Common Reasoning	Open-book-qa	500	Mihaylov et al. (2018)
			Hellaswag	10,003	Zellers et al. (2019)
		Winogrande	1,767	Sakaguchi et al. (2019)	
		PIQA	3,000	Bisk et al. (2020)	
	Finance	Exam	CFA-Easy	1,030	Link
			<b>CFA-Challenge</b>	90	-
<i>Total</i>				91,872	

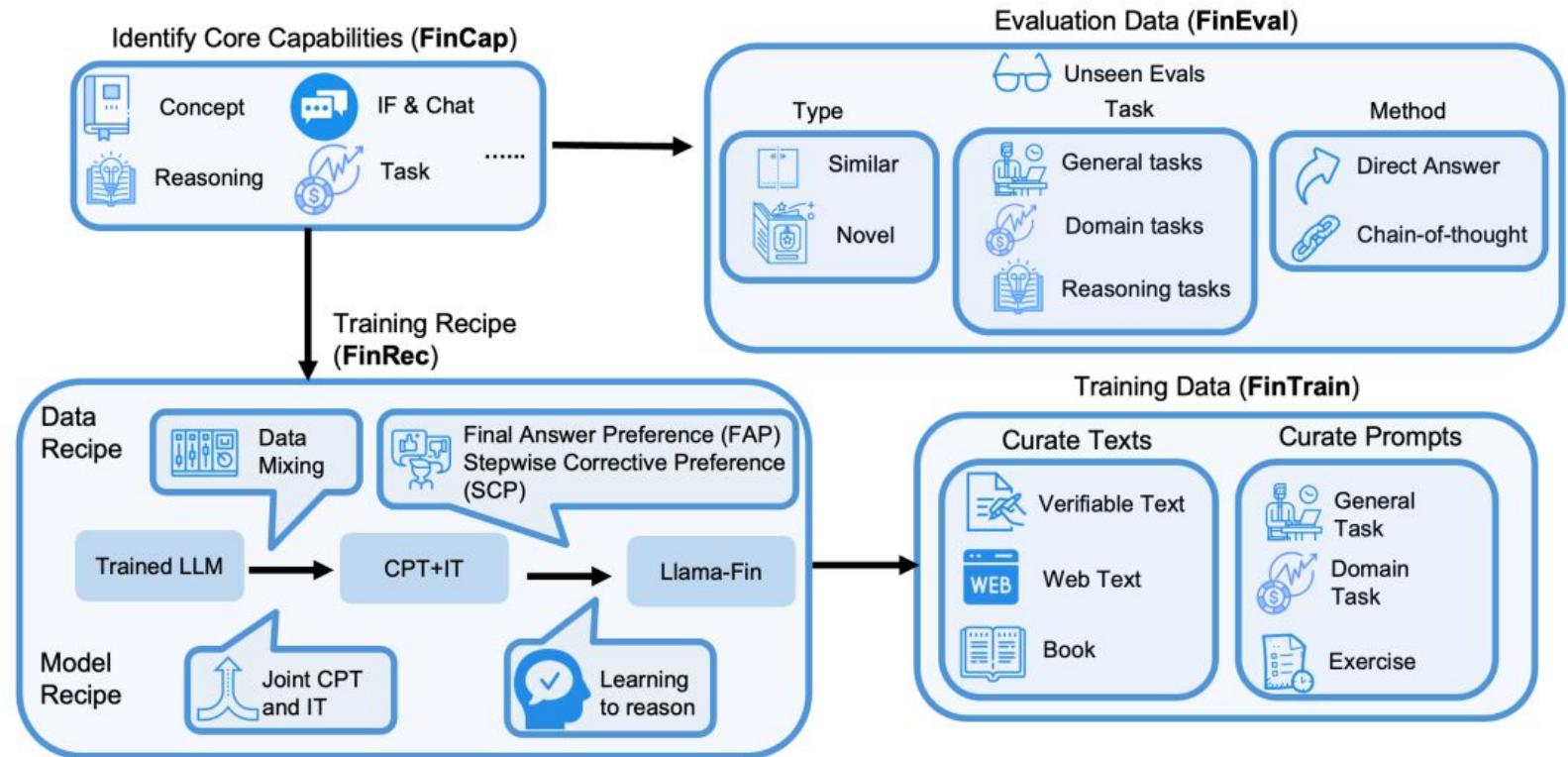
Table 1: Summary of our evaluation dataset. New datasets released with FINDAP are color-highlighted for emphasis.

# FinDAP

## Our Framework



- **FinCap:** Core capacities required for finance domain
- **FinRec:** Our training Recipe
- **FinTrain:** a curated set of training datasets implement FinRec
- **FinEval:** A comprehensive evaluation framework



Finance LLM	Capabilities	Recipe		Evaluation
		Model Recipe	Data Recipe	
AdaptLLM	Concept	CPT	<b>CPT:</b> Financial text + heuristic QAs constructed from the text	Financial + Classification tasks + Direct answer
PIXIU	Task	IT	<b>IT:</b> Financial tasks	Financial + Classification tasks + Direct answer
FinLLM	Concept, Task	CPT → IT	<b>CPT:</b> Financial text + Fineweb; <b>IT:</b> Filtered Financial tasks	Financial + Classification tasks + Direct answer
FinTral	Concept, Task	CPT → IT → PA	<b>CPT:</b> Financial text; <b>IT:</b> Financial tasks; <b>PA:</b> Outcome signal only	Financial + Classification tasks + Direct answer
Palmyra-Fin	SoTA public checkpoint, but recipe is not disclosed			
Llama-Fin	Concept, IF/Chat, Task, Reasoning	CPT+IT → PA	<b>CPT:</b> Financial + General text. <b>IT:</b> Financial + General tasks <b>PA:</b> A novel PA that leverages outcome and process signals	General + Financial tasks; Similar + Novel tasks Classification + Open-form QA tasks Knowledge Recall + Reasoning tasks Direct answer + CoT

Table 1: Comparison between Llama-Fin with other finance LLMs.



# FinDAP

## Final Results (similar tasks)



Outperforms all other baselines (including GPT4o) with one exception

Task	Benchmark	Llama-Fin 8B	Llama3 Instruct 8B	Llama3.1 Instruct 8B	Palmyra Fin 70B	Phi 3.5-mini Instruct 3.8B	Mistral Nemo instruct 12B	GPT4o
Sentiment Ana.	FPB (Acc)	<b><u>91.13</u></b> <sup>✓</sup>	73.09	71.55	67.11	78.04	78.25	82.16
Sentiment Ana.	FiQA SA (Acc)	<b><u>95.32</u></b> <sup>✓</sup>	77.87	70.64	71.91	69.36	55.74	68.51
Monetary Policy	FOMC (Acc)	<b><u>64.31</u></b> <sup>✓</sup>	56.65	54.64	63.10	58.47	57.86	<u>67.94</u>
Named Entity	NER (Rouge1)	<b><u>76.69</u></b> <sup>✓</sup>	45.03	51.22	54.29	39.37	49.84	43.02
Abs Summ.	EDTSUM (Rouge1)	<b><u>53.78</u></b> <sup>✓</sup>	11.50	12.53	21.77	19.97	12.32	18.15

Table 2: Results on **similar (unseen)** tasks. ‘\*’ indicates that ‘GPT4o’ is used as the judge. Llama-Fin and its variant without PA (i.e., the ‘CPT+IT’ checkpoint) are highlighted in blue while the closed model is highlighted in gray . The best performing model for 8b on each benchmark is **bolded**. The overall best performance across all models is underlined. <sup>✓</sup> indicates that Llama-Fin outperforms the base Llama3-8b-inst.



# FinDAP

## Final Results (novel tasks)

General concepts are preserved

Effective in the majority of Tasks  
(12/17)

Instruction-following is also preserved

Excels in reasoning tasks

Capability	Domain	Task	Benchmark	Llama-Fin 8B	Llama3 Instruct 8B	Llama3.1 Instruct 8B	Palmyra Fin 70B	Phi 3.5-mini Instruct 3.8B	Mistral Nemo instruct 12B	GPT4o	
<b>Concept</b>	General	Knowledge Recall	MMLU (CoT, Acc)	47.42	<b>48.14</b>	47.42	54.93	45.07	49.64	<u>63.88</u>	
			AI2-ARC (CoT, Acc)	89.43 <sup>✓</sup>	89.29	<b>89.80</b>	89.01	87.25	88.19	<u>97.85</u>	
			Nq-open (CoT, Acc)	19.20 <sup>✓</sup>	18.47	<b>22.52</b>	19.25	6.20	17.01	<u>27.92</u>	
	Finance	Knowledge Recall	MMLU-Finance (Acc)	64.20	65.71	<b>66.74</b>	75.15	68.17	61.88	<u>86.52</u>	
<b>Task</b>	Finance	Extractive Summ.	Flare-ECTSUM (Rouge1)	34.10	<b>35.92</b>	35.77	33.24	35.52	<u>37.86</u>	35.90	
			ESG Issue	MLESg (Acc)	<b>40.67</b> <sup>✓</sup>	36.33	36.00	39.67	38.33	32.67	<u>45.67</u>
			Rumor Detection	MA (Acc)	84.00 <sup>✓</sup>	82.60	<b>84.20</b>	62.60	75.40	<u>85.20</u>	73.80
		Stock Movement	SM-Bigdata (CoT, Acc)	54.14	<b>55.3</b>	46.06	48.70	53.26	53.53	49.18	
			SM-ACL (CoT, Acc)	<b>51.99</b> <sup>✓</sup>	50.51	45.30	51.21	49.84	50.75	50.97	
			SM-CIKM (CoT, Acc)	54.94	<b>55.56</b>	48.03	52.92	50.03	53.28	49.78	
			Fraud Detection	CRA-CCF (CoT, Mcc)	0.83 <sup>✓</sup>	-0.32	<b>2.73</b>	3.12	1.20	3.94	6.16
		Credit Scoring	CRA-CCFraud (CoT, Acc)	<b>34.03</b> <sup>✓</sup>	14.78	17.3	33.03	45.33	32.94	<u>49.57</u>	
			Flare-German (CoT, Acc)	<b>64.00</b> <sup>✓</sup>	33.50	15.00	12.00	49.50	32.50	17.00	
			Flare-Australian (CoT, Acc)	44.60	<b>66.91</b>	11.51	12.95	46.76	56.12	51.80	
		Distress Ident.	CRA-LendingClub (CoT, Acc)	<b>68.49</b> <sup>✓</sup>	52.69	25.38	23.40	48.87	21.03	65.03	
			CRA-Polish (CoT, Mcc)	<b>15.30</b> <sup>✓</sup>	12.37	15.07	13.78	69.14	11.18	17.38	
			CRA-Taiwan (CoT, Acc)	<b>40.81</b> <sup>✓</sup>	12.01	35.97	52.58	69.96	57.88	8.57	
		Claim Analysis	CRA-ProroSeguro (CoT, Acc)	35.14	<b>96.98</b>	44.33	56.20	25.86	32.58	96.60	
			CRA-TravelInsurance (CoT, Acc)	41.52 <sup>✓</sup>	6.39	80.31	17.28	<b>94.48</b>	73.64	54.03	
		Tabular QA	*Flare-TATQA (CoT, Acc)	<b>66.61</b> <sup>✓</sup>	63.43	63.70	64.21	57.70	66.40	<u>74.90</u>	
			*Finance Bench (CoT, Acc)	<b>54.00</b> <sup>✓</sup>	52.70	38.00	<u>56.67</u>	40.70	55.30	51.30	
<b>IF/Chat</b>	General	Precise IF	MT-bench (1,2 turn avg)	7.36	7.88	<b>7.92</b>	5.80	8.38	7.84	9.10	
<b>Reasoning</b>	Math	Math Reasoning	MathQA (CoT, Acc)	<b>55.08</b> <sup>✓</sup>	51.16	49.35	41.51	39.40	52.46	<u>70.82</u>	
			Social Reasoning	Social-IQA (CoT, Acc)	<b>75.23</b> <sup>✓</sup>	68.83	70.73	77.28	72.82	62.95	<u>78.92</u>
	General	Common Sense	Open-book-qa (CoT, Acc)	<b>82.60</b> <sup>✓</sup>	77.00	82.20	87.00	80.20	76.40	<u>94.60</u>	
			Hellaswag (CoT, Acc)	<b>81.90</b> <sup>✓</sup>	73.34	69.10	69.69	67.89	61.74	81.76	
			Winogrande (CoT, Acc)	<b>70.32</b> <sup>✓</sup>	62.51	66.69	74.27	72.22	65.82	<u>85.71</u>	
			PIQA (CoT, Acc)	<b>85.85</b> <sup>✓</sup>	79.82	81.45	86.72	82.05	77.91	<u>94.34</u>	
	Finance	Exam	CFA-Easy (CoT, Acc)	<b>66.28</b> <sup>✓</sup>	60.56	60.47	36.05	61.24	65.89	<u>83.14</u>	
			CFA-Challenge (CoT, Acc)	<b>55.56</b> <sup>✓</sup>	34.44	35.56	25.56	48.89	43.33	<u>74.44</u>	

Table 3: Results on the **novel** tasks. The notations are the same as in Table 2. ‘Mcc’ refers to Matthews correlation coefficient, usually used in highly imbalanced data (Xie et al., 2024a).

# FinDAP

## Ablations on Preference Alignment

improved on 3 out 5

PA always improve

Mixed: some tasks are inherently 'easy' and reasoning capabilities might not be beneficial (important future work)

PA always improve



Task	Benchmark	Llama-Fin	Llama-Fin (w/o PA)
Sentiment Ana.	FPB	91.13	<b>92.99</b>
Sentiment Ana.	FiQA SA	<b>95.32</b>	94.47
Monetary Policy	FOMC	<b>64.31</b>	63.10
Named Entity	NER	<b>76.69</b>	74.33
Abs. Summ.	EDTSUM	53.78	<b>54.21</b>

Table 4: Ablation on PA on **similar (unseen)** evaluation set.

Capability	Domain	Task	Benchmark	Llama-Fin 8B	Llama-Fin (w/o PA)
<b>Concept</b>	General	Knowledge Recall	MMLU	<b>47.42</b>	47.22
			AI2-ARC	<b>89.43</b>	88.95
			Nq-open	<b>19.20</b>	16.20
	Finance	Knowledge Recall	MMLU-Finance	<b>64.20</b>	63.93
<b>Task</b>	Finance	Extract Summ.	Flare-ECTSUM	34.10	<b>34.41</b>
		ESG Issue	MLESG	40.67	<b>42.00</b>
		Rumor Detection	MA	84.00	<b>84.60</b>
		Stock Movement	SM-Bigdata	<b>54.14</b>	52.04
			SM-ACL	<b>51.99</b>	49.89
			SM-CIKM	<b>54.94</b>	44.88
		Fraud Detection	CRA-CCF	<b>0.83</b>	0.61
			CRA-CCFraud	<b>34.03</b>	32.32
		Credit Scoring	Flare-German	<b>64.00</b>	60.50
			Flare-Australian	44.60	<b>51.80</b>
			CRA-LendingClub	<b>68.49</b>	65.96
		Distress Ident.	CRA-Polish	<b>15.30</b>	0.65
			CRA-Taiwan	40.81	<b>96.41</b>
		Claim Analysis	CRA-ProroSeguro	35.14	<b>86.57</b>
			CRA-TravelInsurance	41.52	<b>98.50</b>
Tabular QA	*Flare-TATQA	<b>66.61</b>	66.43		
Open QA	*Finance Bench	<b>54.00</b>	52.00		
<b>IF/Chat</b>	General	Precise IF	MT-bench	<b>7.36</b>	7.29
<b>Reasoning</b>	Math	Math Reasoning	MathQA	<b>55.08</b>	54.30
		Social Reasoning	Social-IQA	<b>75.23</b>	73.64
	General	Common Sense	Open-book-qa	<b>82.60</b>	79.20
			Hellaswag	<b>81.90</b>	78.92
			Winogrande	<b>70.32</b>	67.48
			PIQA	<b>85.85</b>	84.39
	Finance	Exam	CFA-Easy	<b>66.28</b>	62.31
			CFA-Challenge	<b>55.56</b>	35.56

Table 5: Ablation on PA on **novel** evaluation set.



# FinDAP

## More details

- arXiv (including detailed ablation, hyper-parameters, prompts...):
  - **Demystifying Domain-adaptive Post-training for Financial LLMs**
  - <https://arxiv.org/abs/2501.04961>
- Github: <https://github.com/SalesforceAIResearch/FinDAP>
- HF (FinEval): <https://huggingface.co/datasets/Salesforce/FinEval>

