

# **Bridging On-Device and Cloud LLMs for Collaborative Reasoning**

**Wenzhi Fang, Dong-Jun Han, Liangqi Yuan, Christopher G. Brinton**

# Motivation I: On-Device LLM

- Challenges of Cloud-based LLMs
  - **Pressure on Cloud:** substantial computational pressure on cloud infrastructure
  - **High Latency:** Real-time applications suffer from delays due to network dependency
  - **Underused Device Power:** The potential of device is underexplored
- Advantages of on-device LLMs
  - Offline Capability
  - Better Personalization
  - Enhancing Privacy

# Motivation II: Performance Gap

- The trade-off
  - On-device LLMs
    - Offer privacy, low latency, and device customization
    - But **suffer from limited capacity and weaker performance**
  - Cloud LLMs
    - Deliver **superior accuracy and reasoning power**
    - But introduce latency, privacy risks, and network dependency
- The Key Challenge
  - How to **balance efficiency and accuracy** by leveraging **both on-device and cloud LLMs**

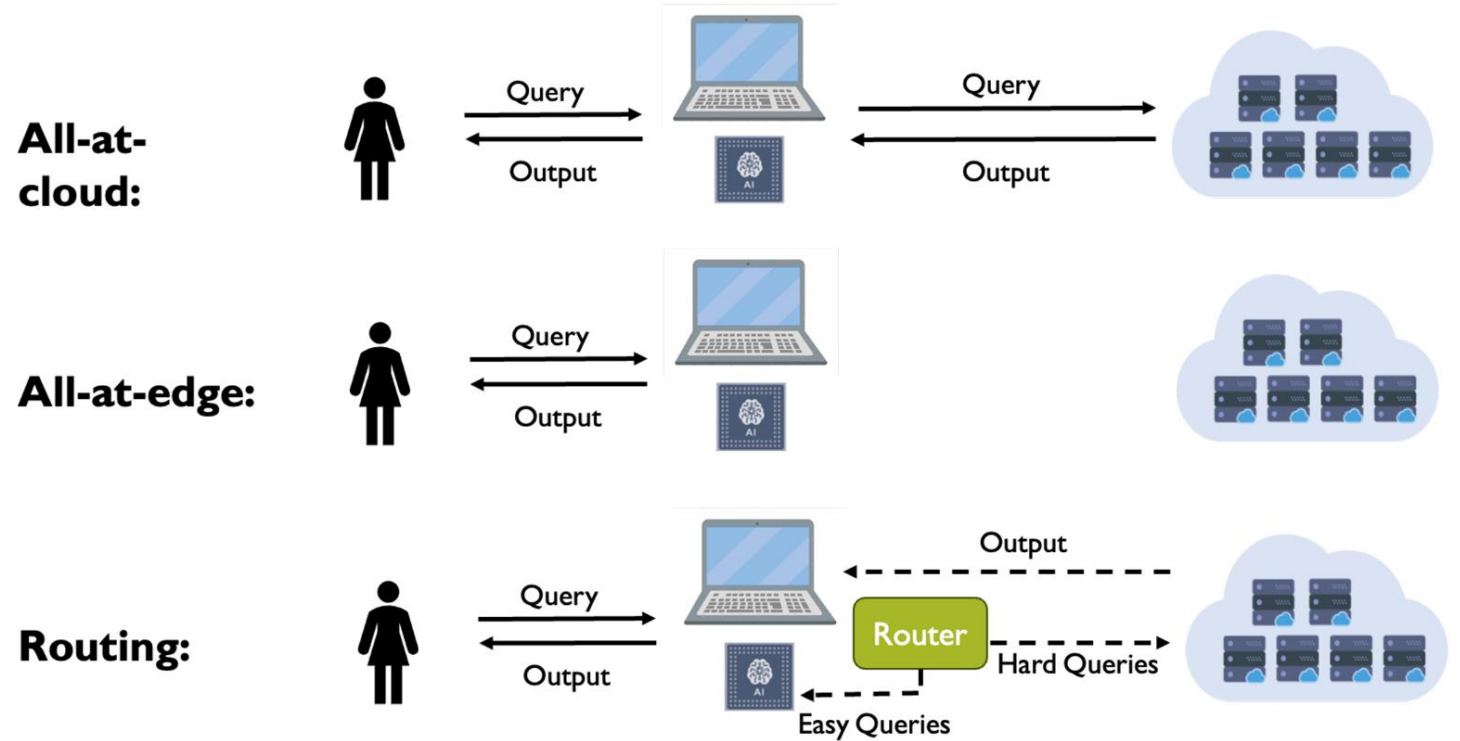
# Existing Pipeline

- Task Optimization

- Prompt tuning
- SFT
- RLHF
- Etc..

- Router optimization

- Binary classifier



# Unified Problem Formulation

- **Key Idea:** The **on-device model** autonomously decide whether to solve the problem independently or **invoke the cloud model**



Using the numbers [74, 78, 36, 7], create an equation that equals 33. You can use basic arithmetic operations (+, -, \*, /) and each number can only be used once.

<think> We want to... Let's try starting with...  
Therefore, the final expression is: </think>

<answer> (74 + 78) - (36 + 7) </answer>



<think> We want to... Start by checking... it seems infeasible... Finally, I am stuck... </think>

<unknown> I need external assistance </unknown>



Using the numbers [74, 78, 36, 7]...



.....  
.....  
.....



**Two possible cases:**

$$y = \begin{cases} y^\theta, & \text{(local only)} \\ [y^\theta, y^c], & \text{(collaborative)} \end{cases}$$

# Unified Problem Formulation

$$\begin{aligned} \max_{\theta} \quad & \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [R(\boldsymbol{\theta}, \mathbf{x})] := \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \mathbb{E}_{\mathbf{y}^{\theta} \sim \pi_{\theta}(\mathbf{x})} [r(\mathbf{x}, \mathbf{y})] \\ \text{s.t.} \quad & \mathbb{E} [\mathbf{1} \{ \mathbf{y} \sim (\pi_{\theta}, \pi_c) \}] \leq \rho \mathbb{E} [\mathbf{1} \{ \mathbf{y} \sim \pi_{\theta} \}]. \end{aligned}$$



Lagrangian Relaxation

Unconstrained Optimization  $\max_{\theta} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \mathbb{E}_{\mathbf{y}^{\theta} \sim \pi_{\theta}(\mathbf{x})} [r(\mathbf{x}, \mathbf{y})] - \gamma \left( \mathbb{E} [\mathbf{1} \{ \mathbf{y} \sim (\pi_{\theta}, \pi_c) \}] - \frac{\rho}{1 + \rho} |\mathcal{D}| \right)$

$$\max_{\theta} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[ \mathbb{E}_{\mathbf{y}^{\theta} \sim \pi_{\theta}(\mathbf{x})} [r(\mathbf{x}, \mathbf{y}) - \gamma \mathbf{1} \{ \mathbf{y} \sim (\pi_{\theta}, \pi_c) \}] \right]$$

$$\max_{\theta} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[ \mathbb{E}_{\mathbf{y}^{\theta} \sim \pi_{\theta}(\mathbf{x})} [\tilde{r}(\mathbf{x}, \mathbf{y})] \right] \quad \tilde{r}(\mathbf{x}, \mathbf{y}) = \begin{cases} r(\mathbf{x}, \mathbf{y}) - \gamma, & \text{if } \mathbf{y} \sim (\pi_{\theta}, \pi_c), \\ r(\mathbf{x}, \mathbf{y}), & \text{otherwise.} \end{cases}$$

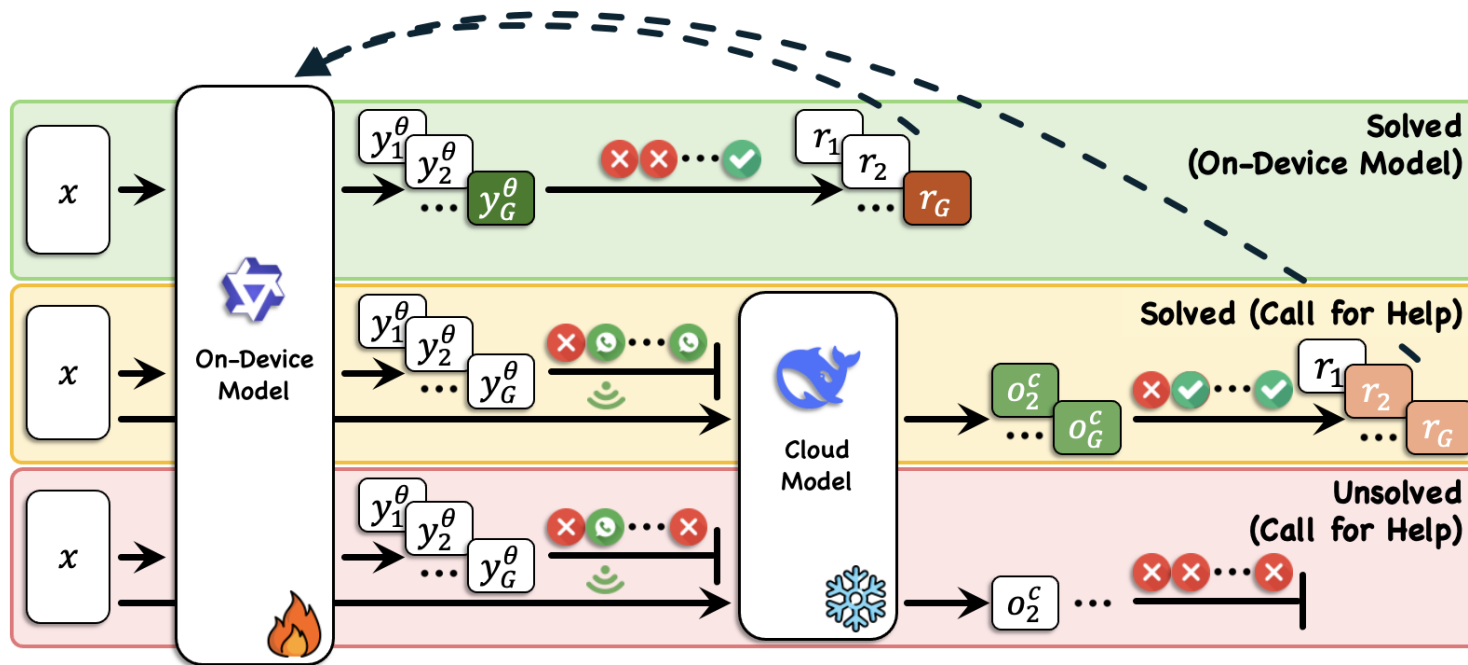
# Unified Problem Formulation

$$\max_{\theta} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[ \mathbb{E}_{\mathbf{y}^{\theta} \sim \pi_{\theta}(\mathbf{x})} [\tilde{r}(\mathbf{x}, \mathbf{y})] \right] \quad \tilde{r}(\mathbf{x}, \mathbf{y}) = \begin{cases} r(\mathbf{x}, \mathbf{y}) - \gamma, & \text{if } \mathbf{y} \sim (\pi_{\theta}, \pi_c), \\ r(\mathbf{x}, \mathbf{y}), & \text{otherwise.} \end{cases}$$

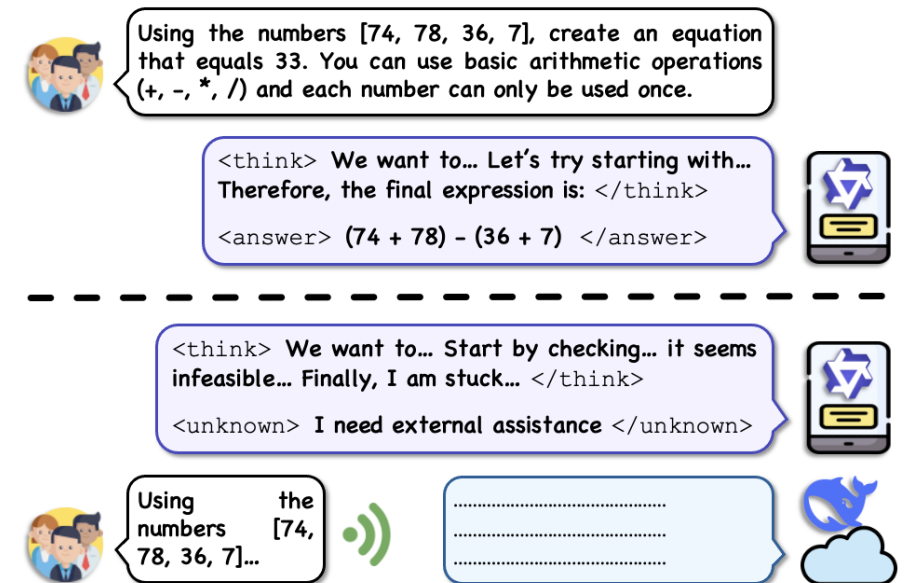
- Two possible case  $y = \begin{cases} y^{\theta}, & \text{(local only)} \\ [y^{\theta}, y^c], & \text{(collaborative)} \end{cases}$
- **Accuracy Reward (the largest)**: Rewards the correctness of the on-device model's final answer
- **Coordination Reward**: Rewards effective delegation, granted when the device invokes the cloud model and the joint output is correct

# Collaborative Unified Training Framework

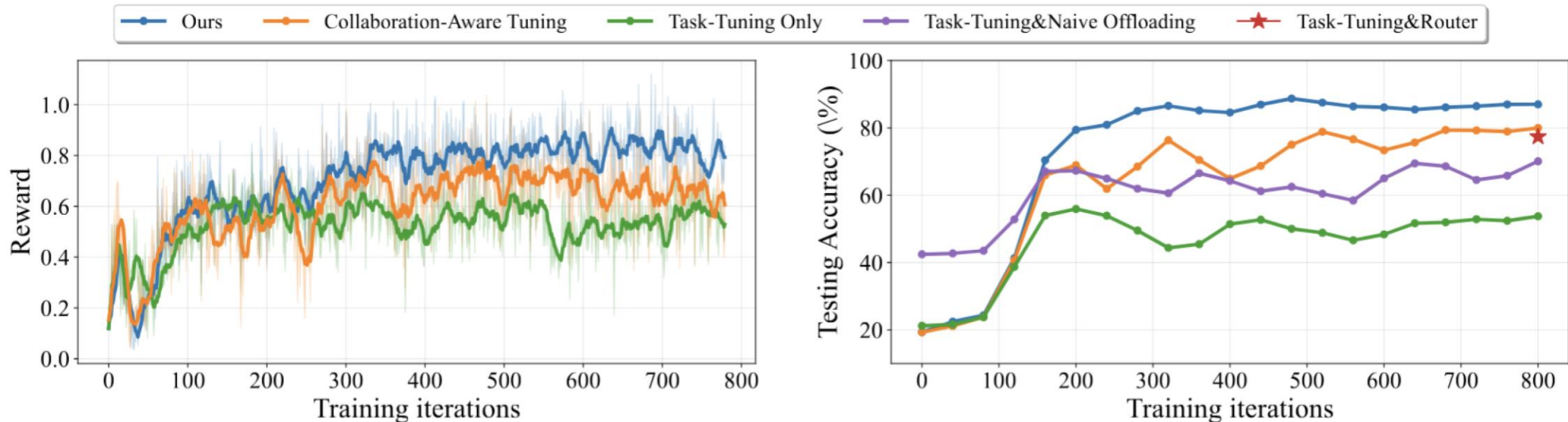
- Our methodology



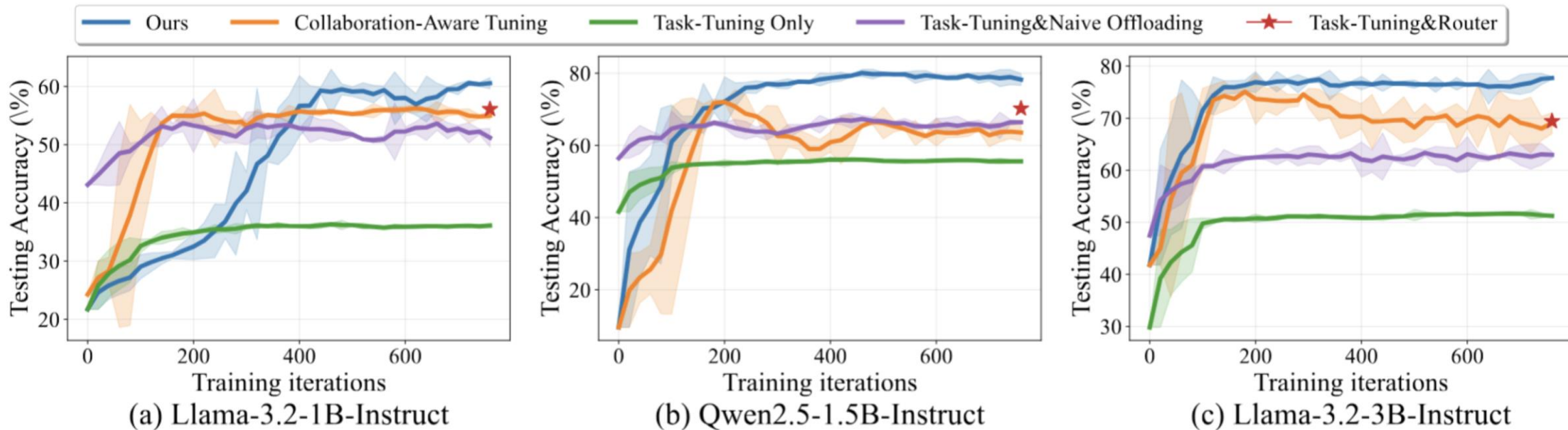
(a) Our Training Methodology



(b) Collaborative Inference



### Qwen2.5-3B-Instruct



Model	Metric	Method	MATH-lighteval	MATH-500	AMC23	MinervaMath	AGI-Eval-Math	Avg.
		Cloud LLM	98.4	97.3	97.5	80.9	94.7	93.8
Qwen2.5-1.5B	Cond. Acc.	Task-Tuning Only	56.1	54.8	35.0	20.6	51.8	43.7
		Task-Tuning&Naive Offloading	56.3	55.4	32.1	21.0	51.6	43.3
		Collaboration-Aware Tuning	54.1	61.2	28.6	13.6	55.3	42.6
		Task-Tuning&Router	61.8	64.9	39.3	21.8	58.7	49.3
		<b>Ours</b>	<b>72.6</b>	<b>75.1</b>	<b>42.9</b>	<b>24.1</b>	<b>64.6</b>	<b>55.8</b>
	Overall Acc.	Task-Tuning Only	56.1	54.8	35.0	20.6	51.8	43.7
		Task-Tuning&Naive Offloading	67.2	67.4	50.0	38.2	68.3	58.2
		Collaboration-Aware Tuning	61.5	61.2	42.5	33.5	66.9	53.1
		Task-Tuning&Router	70.9	72.2	55.0	36.8	69.2	60.8
		<b>Ours</b>	<b>80.4</b>	<b>81.6</b>	<b>57.5</b>	<b>40.8</b>	<b>73.4</b>	<b>66.7</b>
Llama-3.2-3B	Cond. Acc.	Task-Tuning Only	51.2	43.0	27.5	19.1	45.5	37.3
		Task-Tuning&Naive Offloading	51.6	41.8	27.5	19.4	45.7	37.2
		Collaboration-Aware Tuning	60.2	43.7	21.4	13.4	43.9	36.5
		Task-Tuning&Router	64.9	45.7	25.0	20.2	35.1	38.2
		<b>Ours</b>	<b>72.2</b>	<b>56.6</b>	<b>35.7</b>	<b>27.7</b>	<b>59.7</b>	<b>50.4</b>
	Overall Acc.	Task-Tuning Only	51.2	43.0	27.5	19.1	45.5	37.3
		Task-Tuning&Naive Offloading	65.1	59.0	45.0	37.1	58.7	53.0
		Collaboration-Aware Tuning	66.8	59.6	42.5	36.8	58.9	52.9
		Task-Tuning&Router	69.4	62.0	46.0	39.1	53.8	54.1
		<b>Ours</b>	<b>79.5</b>	<b>68.6</b>	<b>52.5</b>	<b>43.4</b>	<b>64.5</b>	<b>61.7</b>

**The proposed method generalize well to the out-of-distribution data**