國立陽明交通大學
NATIONAL YANG MING CHIAO TUNG UNIVERSITY

**Paper Tag: WED-PM-247**

# Multimodal Prompting with Missing Modalities for Visual Recognition

**Yi-Lun Lee**          **Yi-Hsuan Tsai**          **Wei-Chen Chiu**          **Chen-Yu Lee**
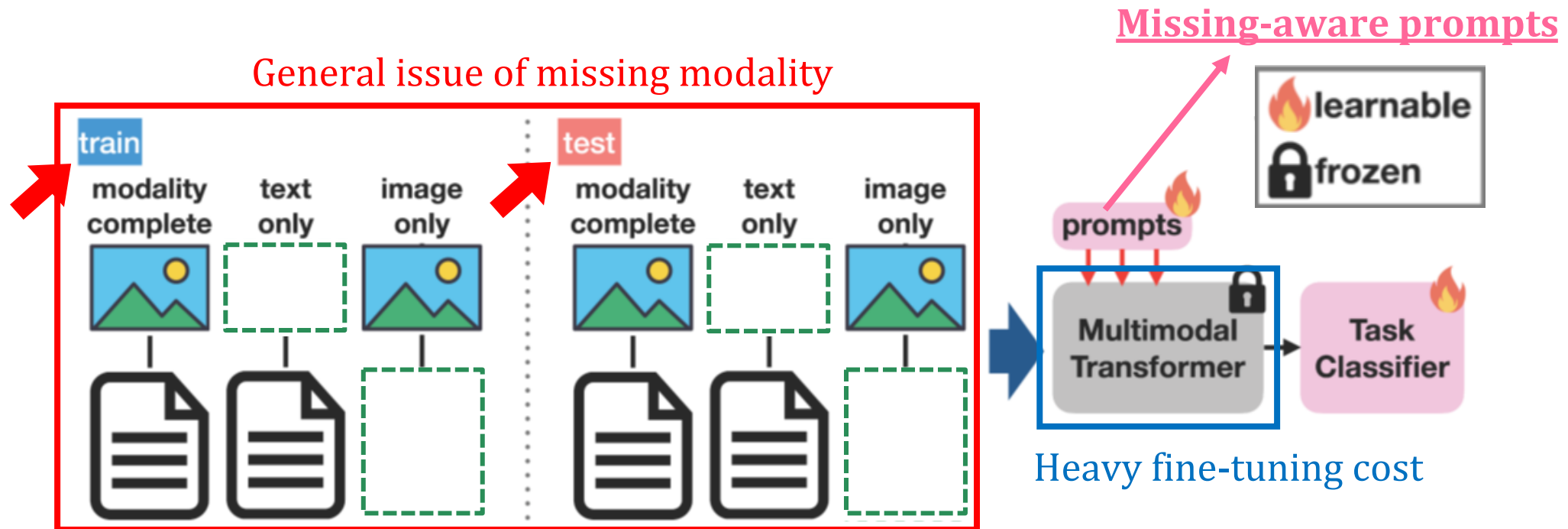
Enriched Vision Applications Laboratory

Google

# Goal

- A simple <u>prompt-learning-based</u> method for multimodal learning:
  - Tackle the general issue of missing modality
  - No need to finetune the heavy pre-trained model (transformer)

**Missing-aware prompts**

General issue of missing modality



Heavy fine-tuning cost

# Multimodal Learning

- Our observation perceived in daily life is typically multimodal

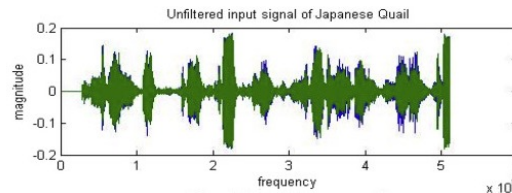| Modality | Data sample | Properties | Target task |
|---|---|---|---|



**Multimodal Visual Recognition**

Appearance

[From Wikipedia] The morphology of the Japanese quail differs depending on its stage in life. As chicks, both male and female individuals exhibit the same kind of plumage and coloring...
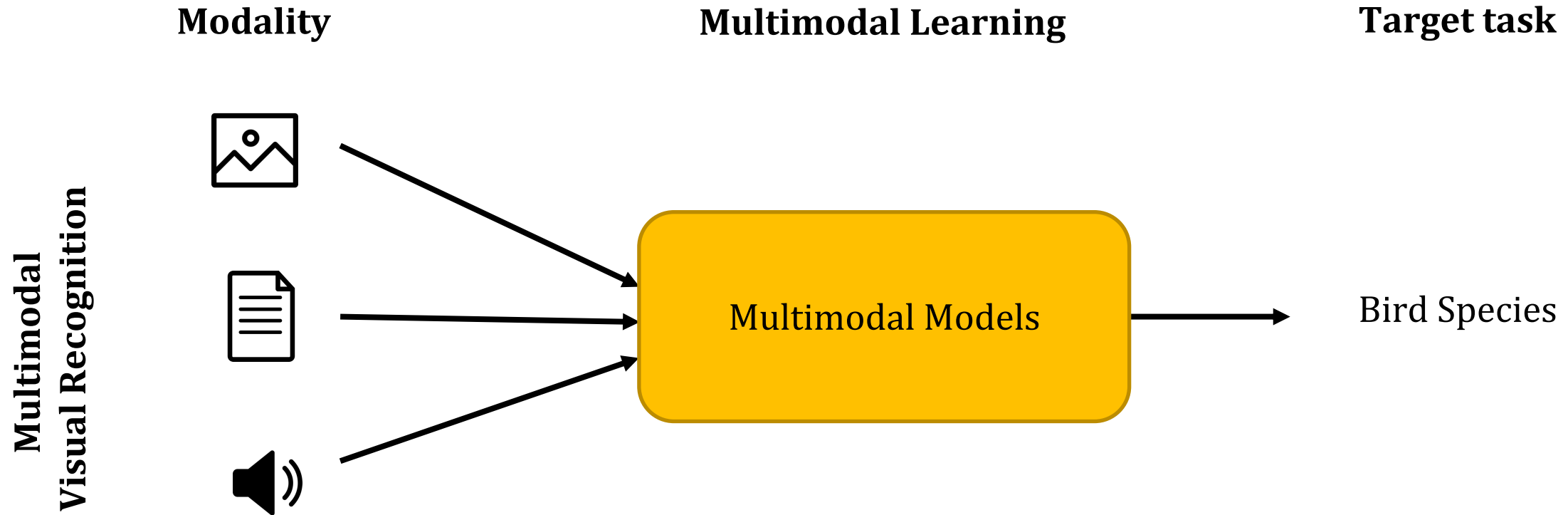
Description

Bird Species

Sound

# Multimodal Learning

- *GOAL:* leverage the potential complementary properties among modalities to better realize the target tasks

**Modality**          **Multimodal Learning**          **Target task**

**Multimodal Visual Recognition**



Multimodal Models

Bird Species

# Challenges

- Some practical challenges for multimodal methods
  - Missing modality
  - Heavy cost of finetuning huge pre-trained models
  - Noisy web-crawled data (i.e. incomplete and incorrect)
  - Multimodal data perturbation (i.e. distribution shifts in real world)

# Challenge: Missing Modality

- Missing modality could happen for different reasons
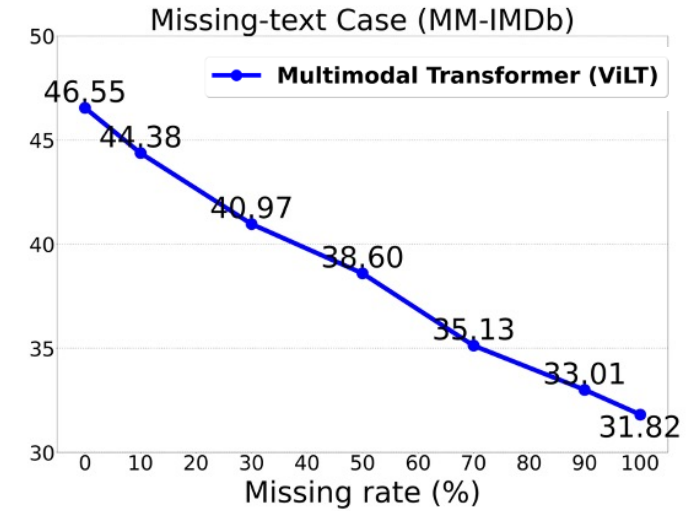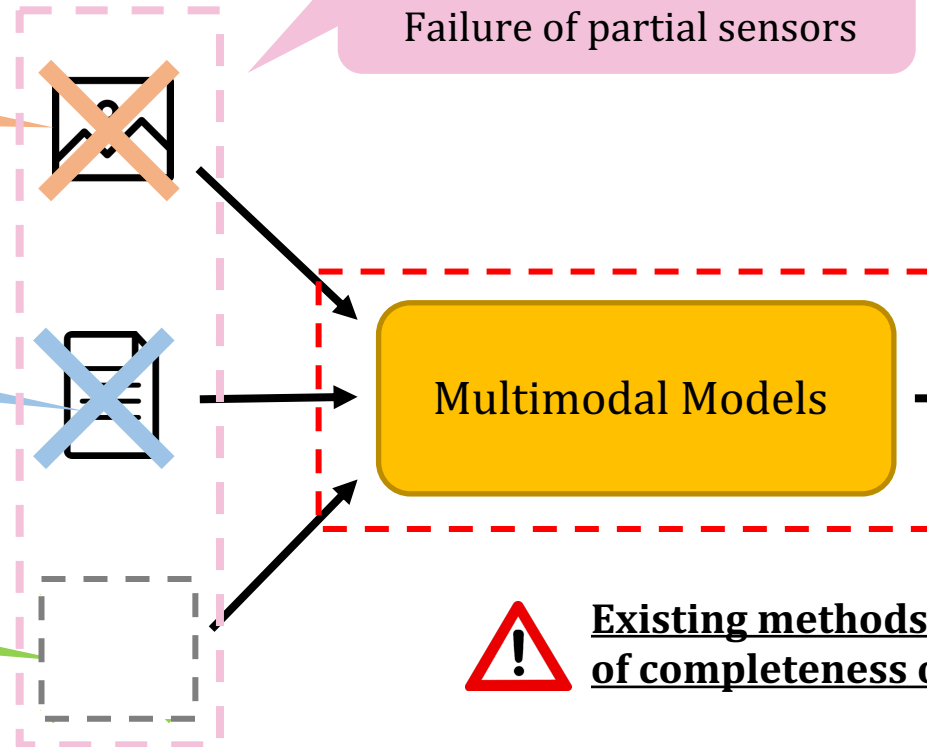
Modality-specific reasons

General reasons

Asynchrony among sensors
Failure of partial sensors

Non-coverage of camera
Occlusion

Text missing
Privacy issue

Multimodal Models

Low sound
Ambient noise

Missing-text Case (MM-IMDb)

Multimodal Transformer (ViLT)

46.55
44.38
40.97
38.60
35.13
33.01
31.82

Missing rate (%)

Performance drop

Fail to predict

⚠ **Existing methods usually have the assumption of completeness of multi-modality.**

# Challenge: Heavy Training Cost

- Heavy cost of finetuning pre-trained models
  - Huge size: billions of parameters
    - E.g. GPT-3 has 175B parameters
  - Long finetuning time
  - Generalization ability (overfitting issue, stability issue)
- If we only have limited computation resource…

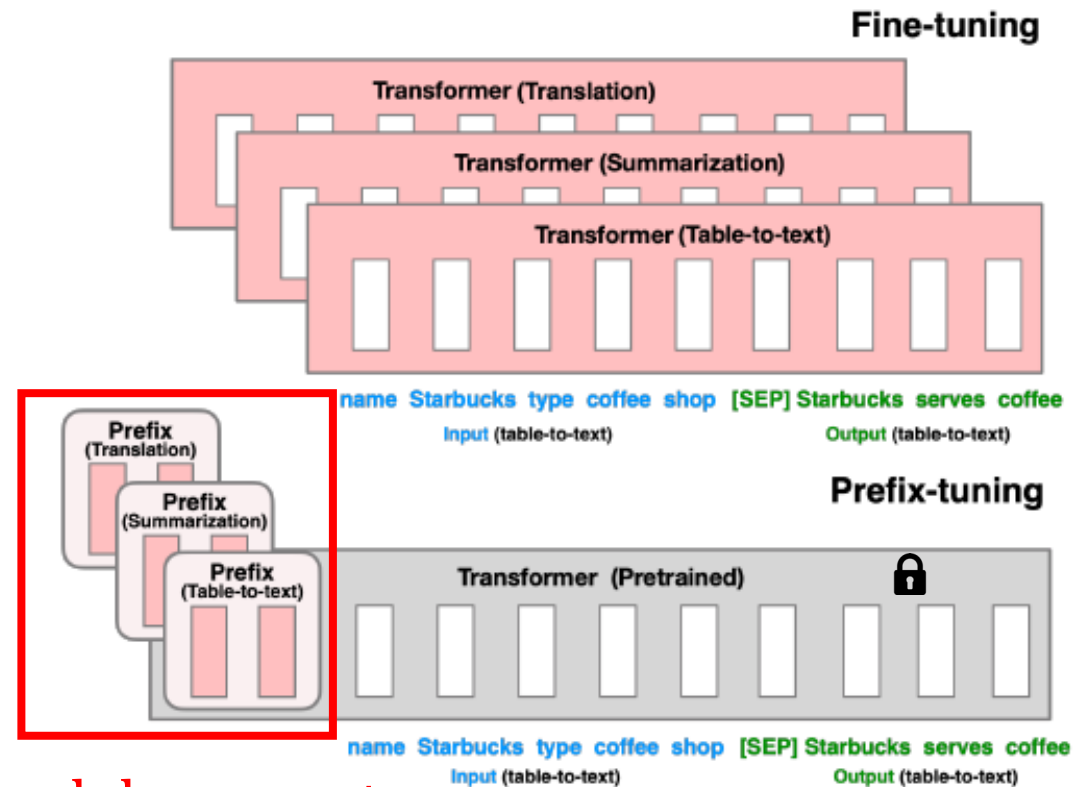Q: How can we efficiently and effectively finetune pre-trained models?
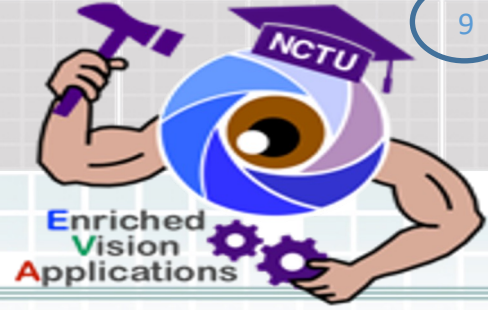
**Prompt learning**

# Prompt Learning

- Learnable "task prompts" instruct models to perform specific downstream tasks



Li et al. "Prefix-Tuning: Optimizing Continuous Prompts for Generation"
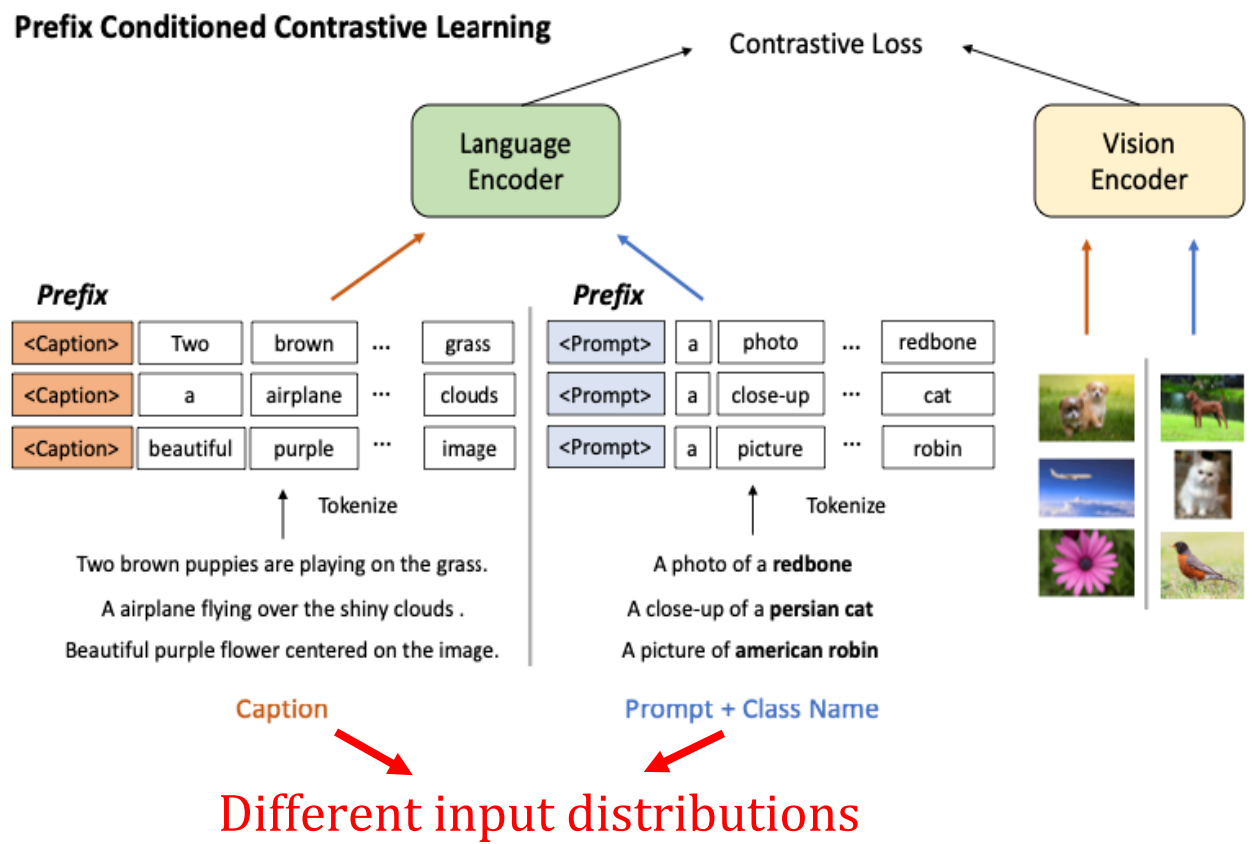
# Prompt Learning

- Different prompts instruct the model learning with <u>different input distributions</u>



**Different input distributions**

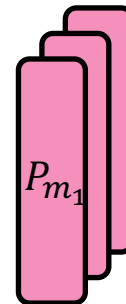Saito et al. "Prefix Conditioning Unifies Language and Label Supervision"

# Motivation

- Missing modalities can be regarded as <u>different input distributions</u>

  - Complete: real text + real image

  - Text-only: real text + dummy image

  - Image-only: dummy text + real image

- Use prompts to learn with modality-missing data accordingly
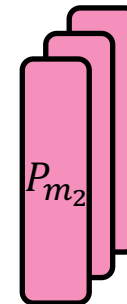


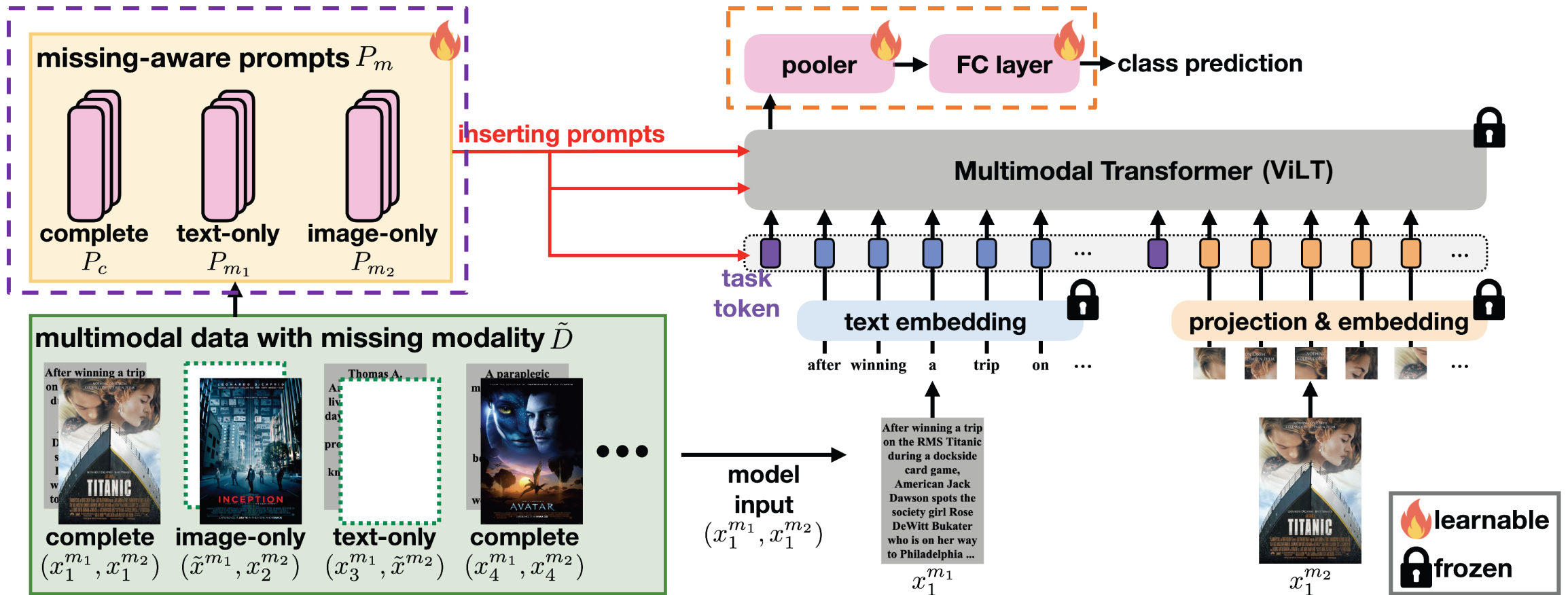Complete        Text-only        Image-only
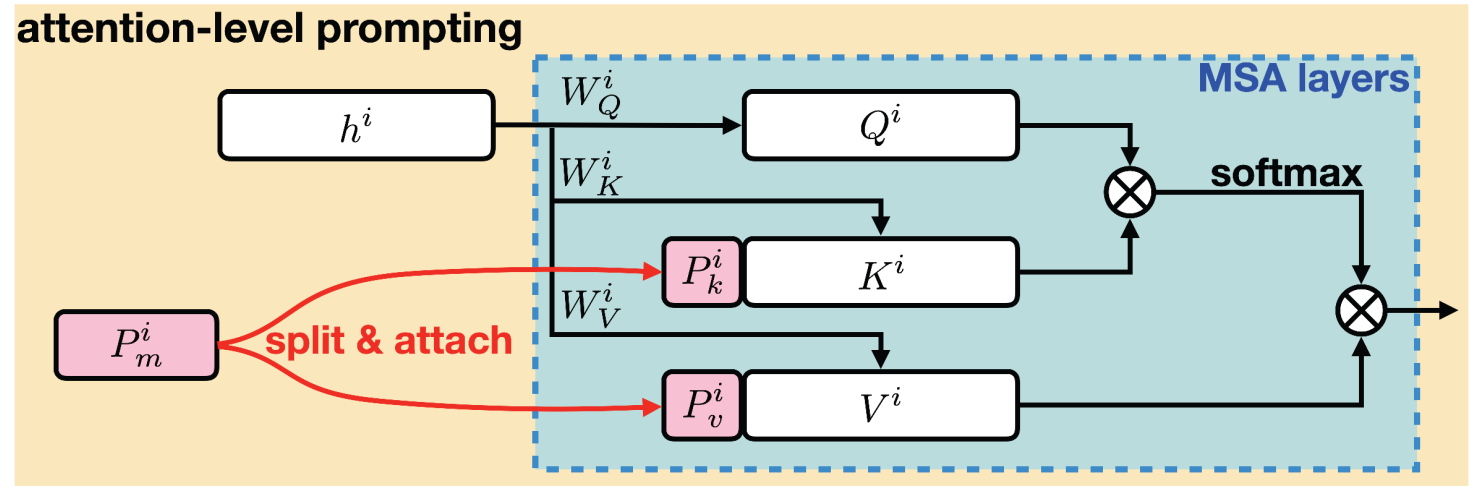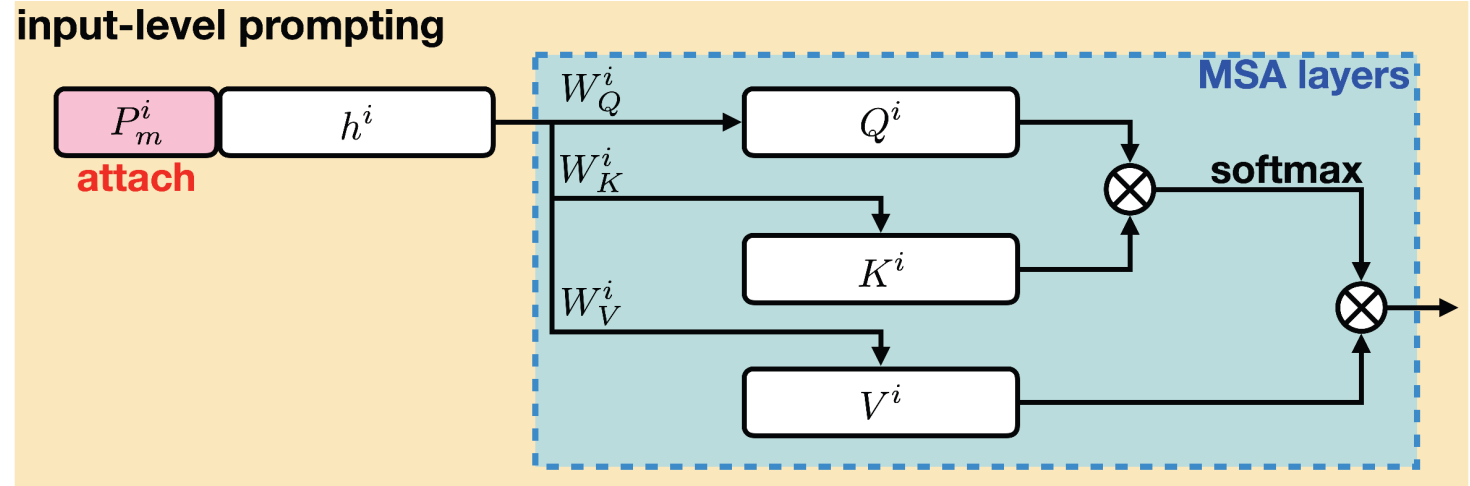
# Proposed Method

- Our prompting framework

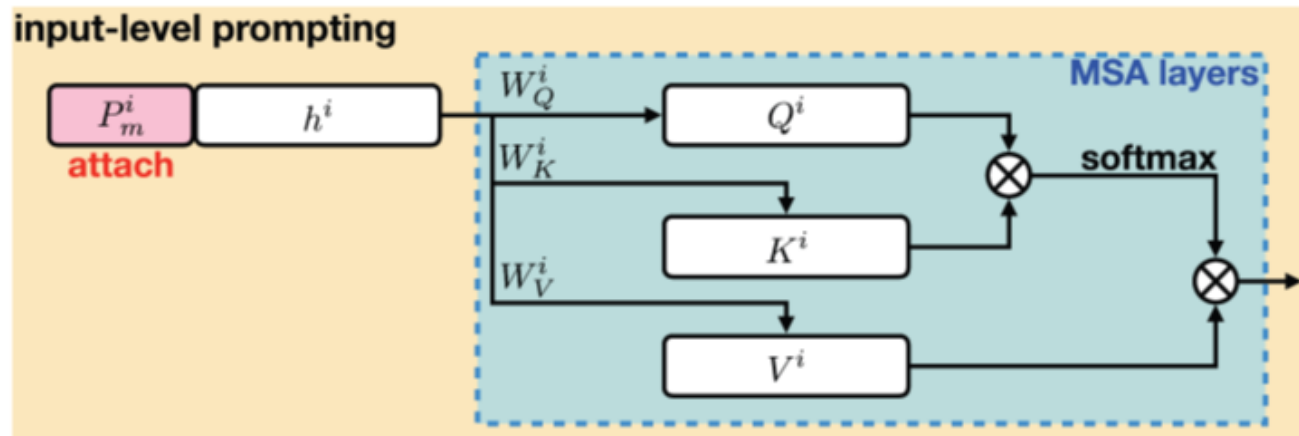$$L = L_{task}(x_i^{m_1}, x_i^{m_2}, \boxed{\theta_t}, \boxed{\theta_p})$$

# Proposed Method

- Prompt design
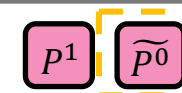  - Input-level
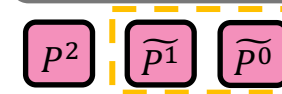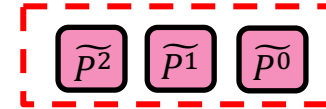  - Attention-level

# Proposed Method

- Input-level prompting
  - Inheriting instruction information from previous layers could be helpful
  - increasing sequence length
    - Sensitive to different datasets
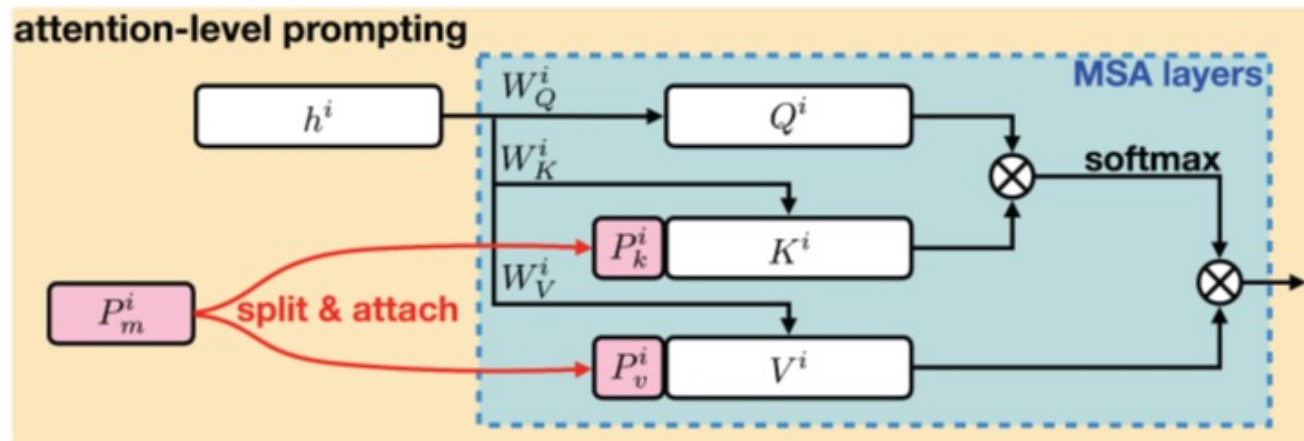
# Proposed Method

- Attention-level prompting
  - Insert the prompts into key and value of MSA layers
  - Focus on current layer instruction
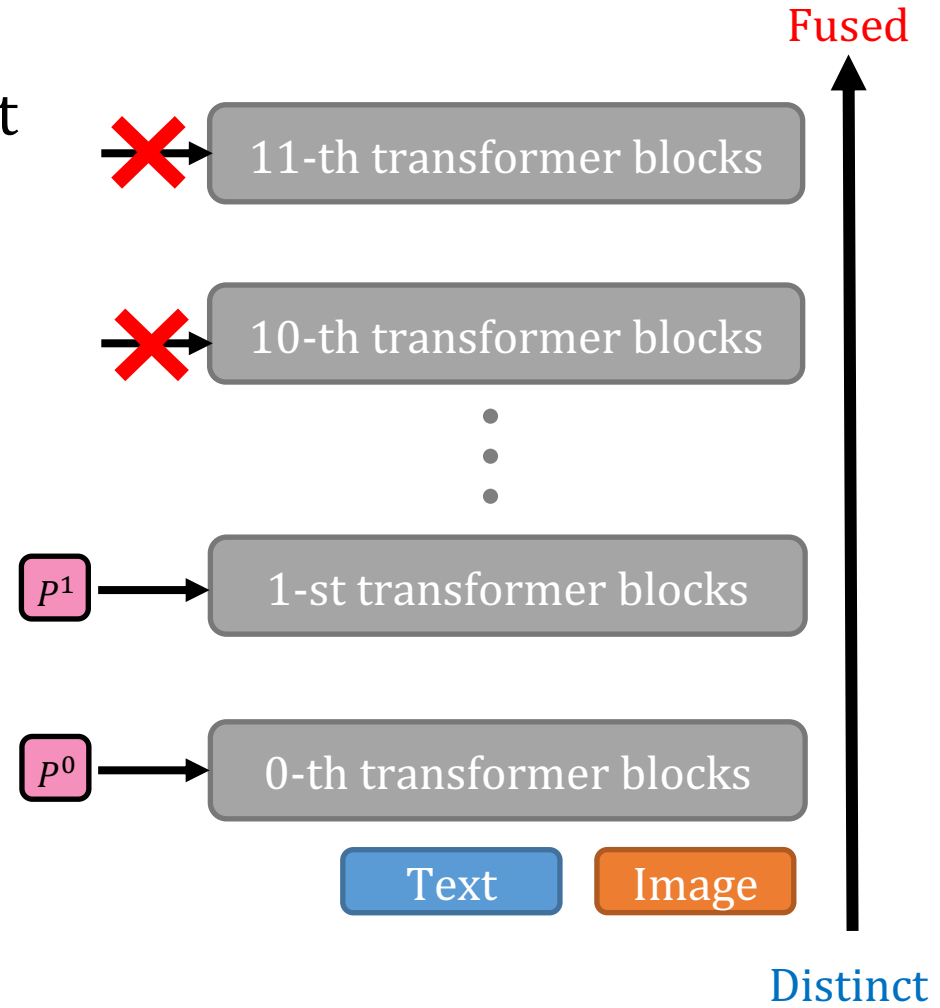  - No increasing length -> less-sensitive to datasets



$$Attention^l = softmax\left(\frac{Q^l[p_k^l; K^l]^T}{\sqrt{d}}\right)[p_v^l; V^l]$$

$L \times D \qquad L \times (L + P) \qquad (L + P) \times D$

# Proposed Method

- Location for multi-layer prompting

  - The features of different layers could be different

    - Earlier layer features are more distinct

    - Later layer features are more well-fused

  - Prompting in the earlier layer is the choice

Fused

| ✖ → | 11-th transformer blocks |

| ✖ → | 10-th transformer blocks |

$P^1$ → | 1-st transformer blocks |

$P^0$ → | 0-th transformer blocks |

| Text | Image |

Distinct

# Experiments

- Multimodal vision recognition datasets
  - MM-IMDb – multi-label classification
  - UPMC Food-101 – single-label classification
  - Hateful Memes – binary classification

| Datasets | Text length | Image length |
|---|---|---|
| MM-IMDb | 1024 | 192-216 |
| UPMC Food-101 | 512 | 192-216 |
| Hateful Memes | 128 | 192-216 |



MM-IMDb



UPMC Food-101



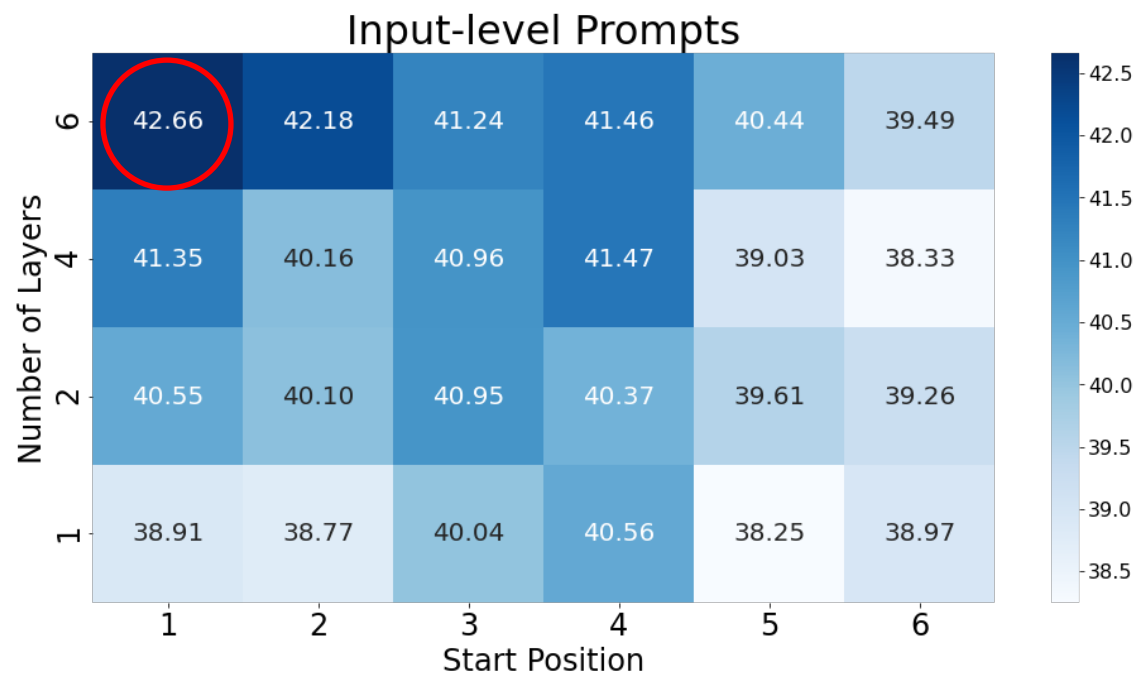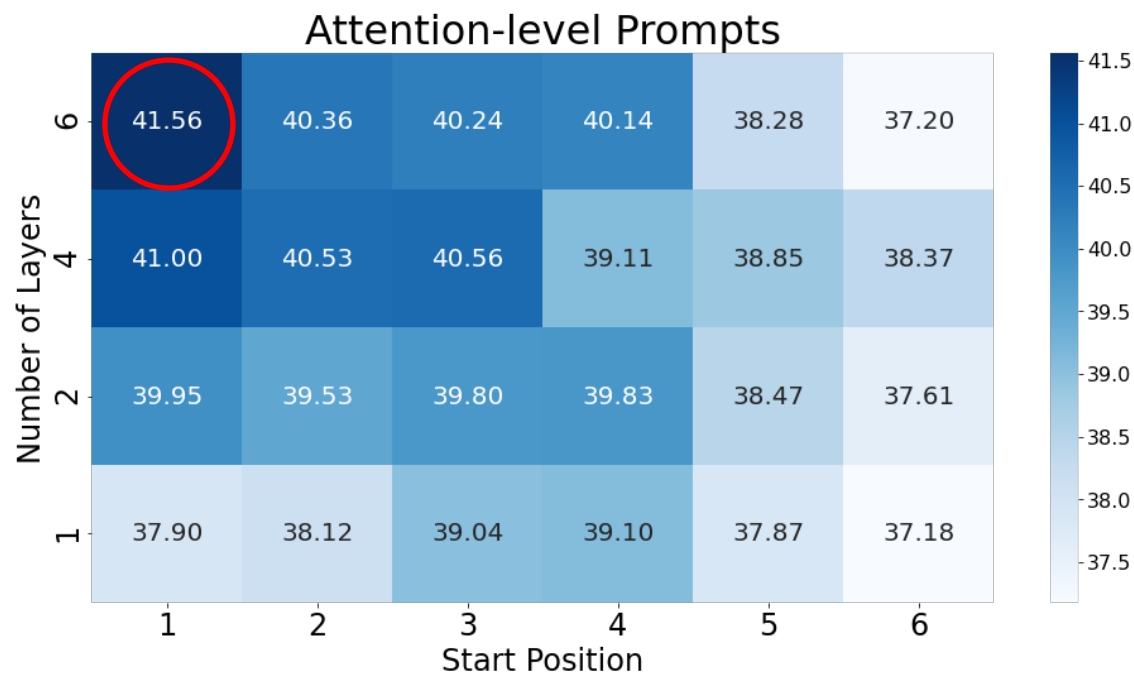Hateful Memes

# Quantitative Results

- Baseline
  - Pre-trained ViLT
  - Only train task-related models (i.e. classifier)

- Input-level prompts
  - Better performance
  - Sensitive to datasets

- Attention-level prompts
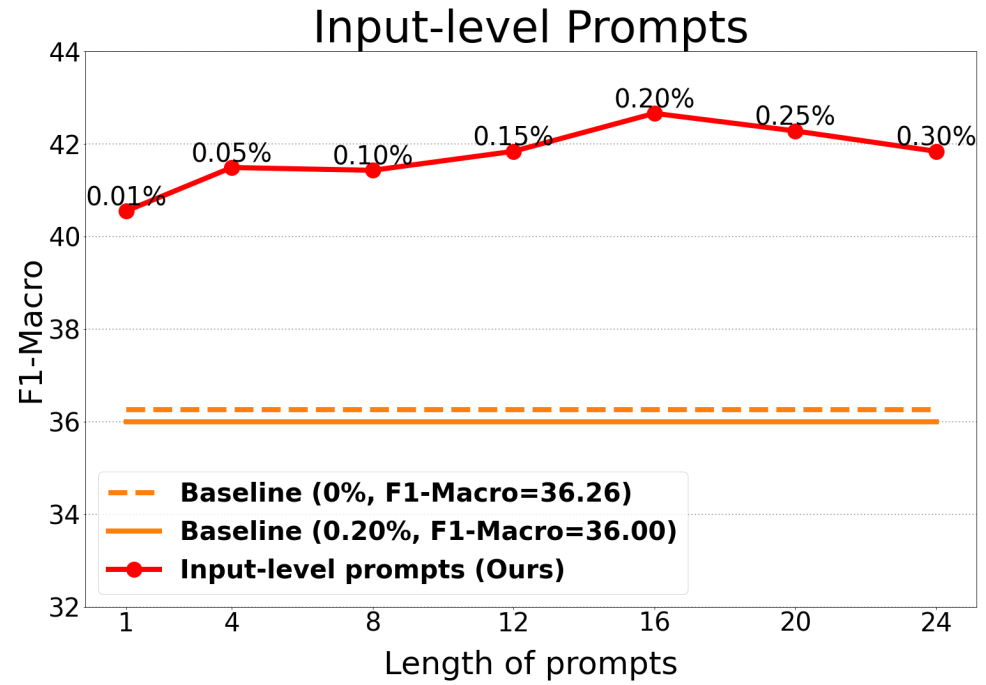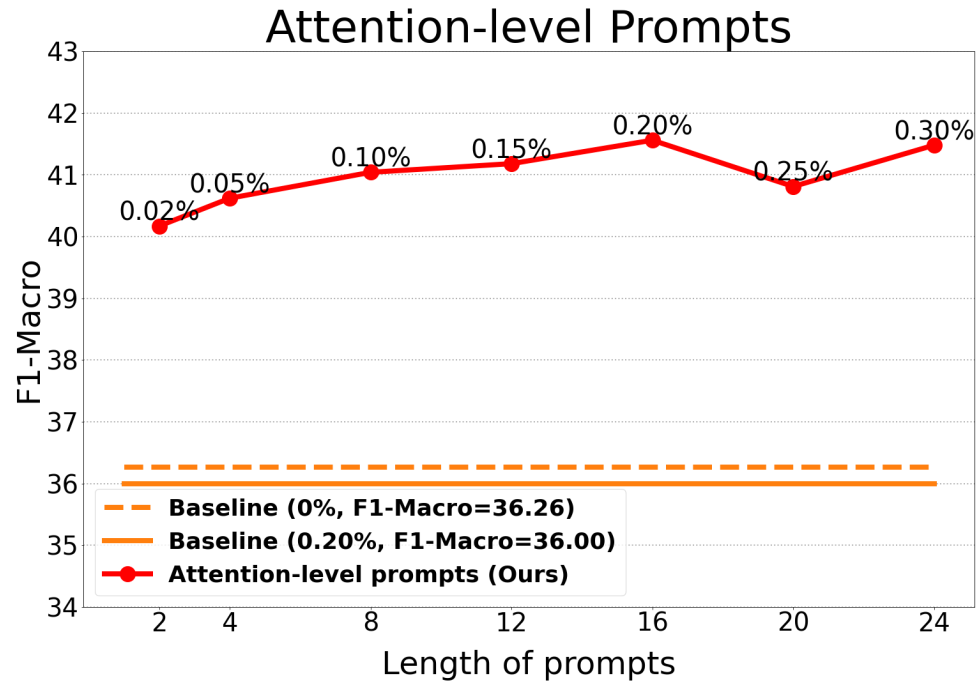  - Consistently improvement

# Ablation Study

- Prompting position
  - the earlier prompting layers and more prompting layers improve the performance



Attention-level Prompts

| Number of Layers \ Start Position | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| 6 | 41.56 | 40.36 | 40.24 | 40.14 | 38.28 | 37.20 |
| 4 | 41.00 | 40.53 | 40.56 | 39.11 | 38.85 | 38.37 |
| 2 | 39.95 | 39.53 | 39.80 | 39.83 | 38.47 | 37.61 |
| 1 | 37.90 | 38.12 | 39.04 | 39.10 | 37.87 | 37.18 |

Input-level Prompts

| Number of Layers \ Start Position | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| 6 | 42.66 | 42.18 | 41.24 | 41.46 | 40.44 | 39.49 |
| 4 | 41.35 | 40.16 | 40.96 | 41.47 | 39.03 | 38.33 |
| 2 | 40.55 | 40.10 | 40.95 | 40.37 | 39.61 | 39.26 |
| 1 | 38.91 | 38.77 | 40.04 | 40.56 | 38.25 | 38.97 |

# Ablation Study

- Prompt length
  - even with <span style="color:red">fewer parameters</span> (i.e., reducing the prompt length to 1), the performance <span style="color:red">is still competitive</span>