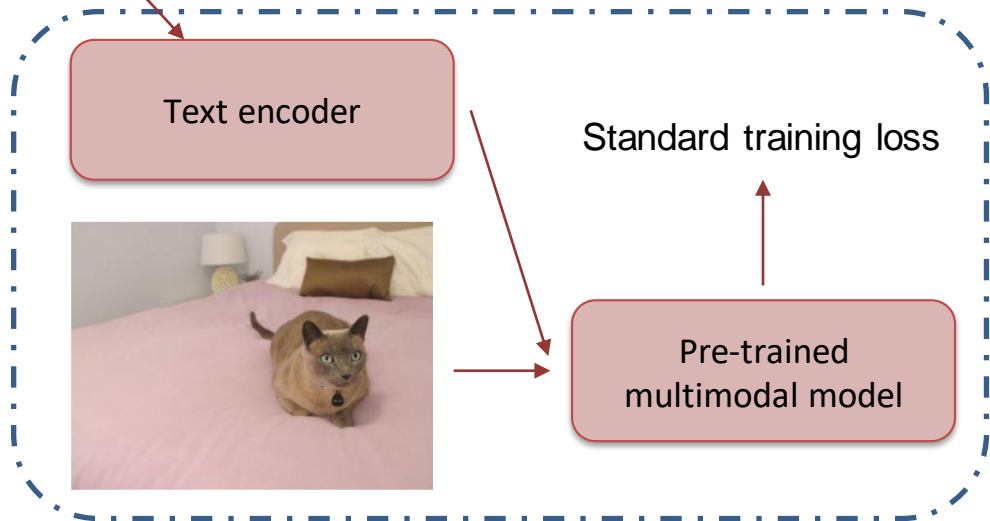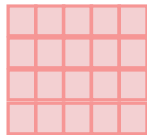JUNE 18-22, 2023
CVPR
VANCOUVER, CANADA

# Learning to Name Classes for Vision and Language Models

Sarah Parisot, Yongxin Yang, Steven McDonagh

Huawei Noah's Ark Lab

[prompt context] + [class name]

Adapt vision-language models to new dataset by **learning class names**

Text encoder

Standard training loss

Pre-trained multimodal model
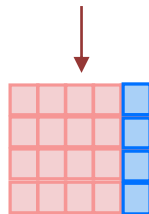
Adapt vision-language models to new dataset by **learning class names**

- Removes class name ambiguities

- Increases robustness to prompt context

- Language agnostic: adapt to model's observed language

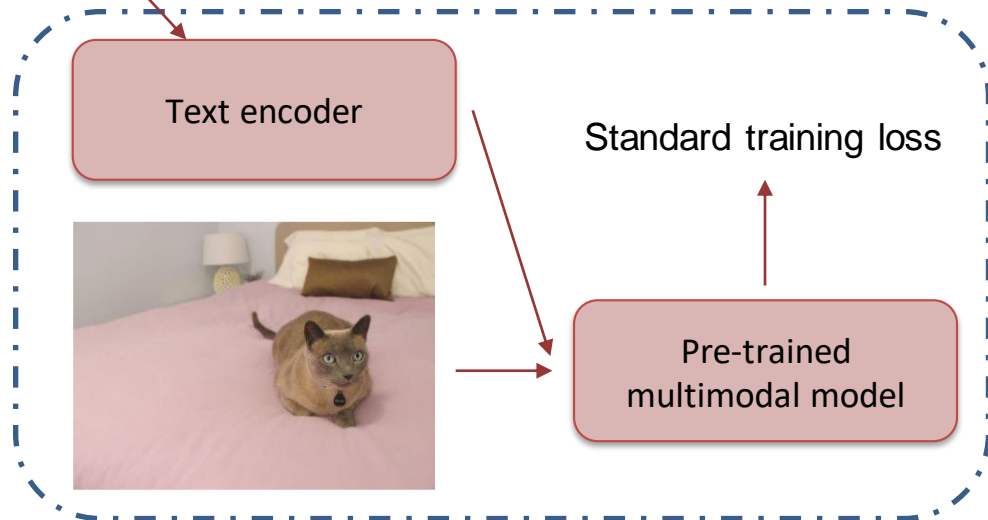- Directly applicable to both classification and object detection tasks
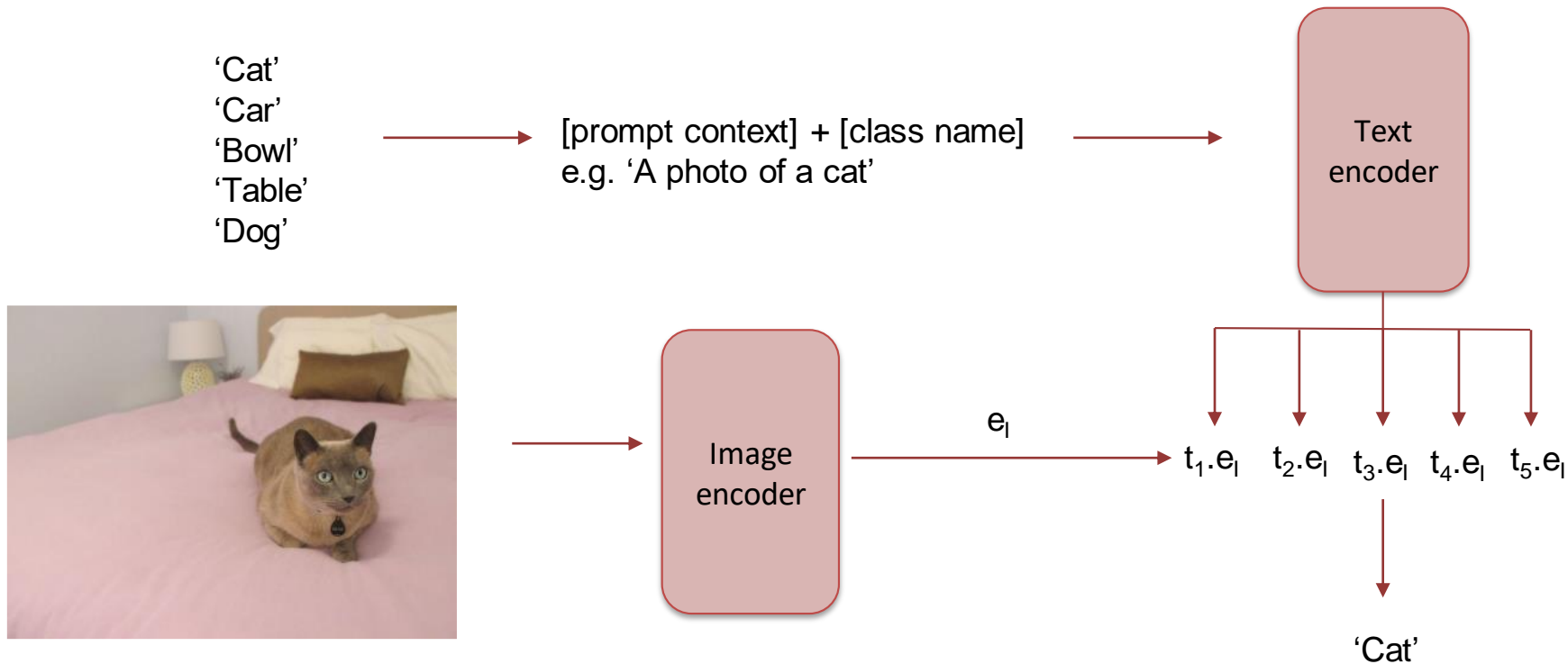
[prompt context] + [Placeholder]
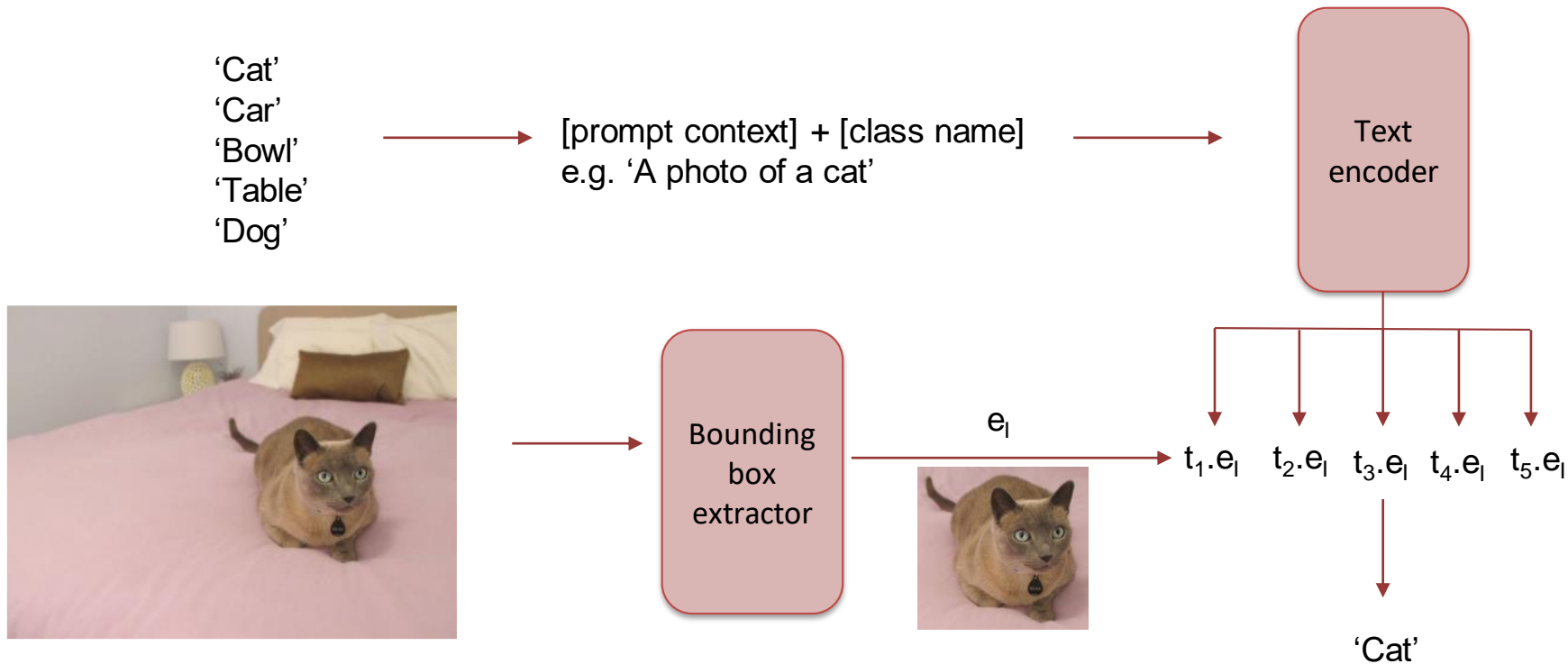
Learnable word embeddings

: Cat
: Car
: Table
: Bed

Text encoder

Standard training loss

Pre-trained multimodal model

# Vision-language classification models

'Cat'
'Car'
'Bowl'
'Table'
'Dog'

[prompt context] + [class name]
e.g. 'A photo of a cat'

Text encoder

Image encoder

$e_I$

$t_1.e_I$  $t_2.e_I$  $t_3.e_I$  $t_4.e_I$  $t_5.e_I$

'Cat'

Radford et al. "Learning transferable visual models from natural language supervision." *International Conference on Machine Learning*. PMLR, 2021.

# Vision-language detection models

‘Cat’
‘Car’
‘Bowl’
‘Table’
‘Dog’

[prompt context] + [class name]
e.g. ‘A photo of a cat’

Text encoder

Bounding box extractor

$e_I$

$t_1.e_I$  $t_2.e_I$  $t_3.e_I$  $t_4.e_I$  $t_5.e_I$

‘Cat’

Minderer et al. "Simple open-vocabulary object detection with vision transformers." *ECCV* 2022
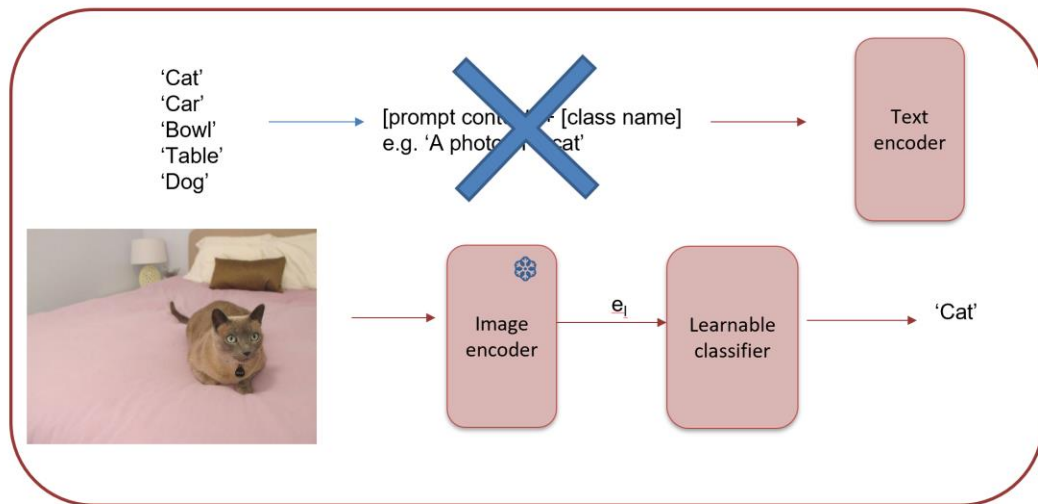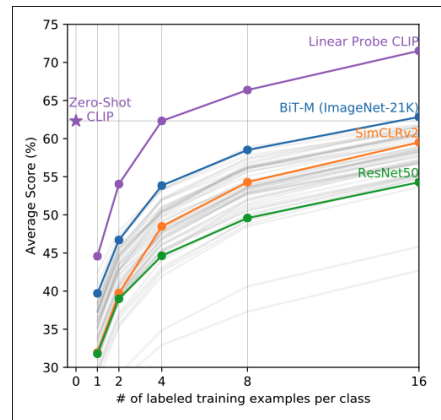
# Fine-tuning



- Adapting vision-language models to new data: **challenging!**
  - Small dataset overfitting
  - Losing generalisation ability

- Linear probing
  - Train a standard linear classification layer using frozen image encoder
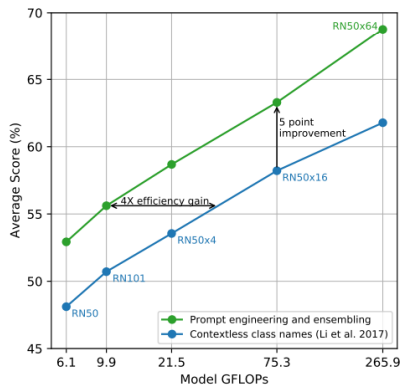
  ✅ Data efficient
  ✅ improves over zero-shot performance
  ✅ no hand-crafted text components

  ❌ Loses open-set and zero-shot properties

# Sensitivity to prompt input

- Model performance is sensitive to text input



- Existing methods rely on *handcrafted* class names. Potentially:
  - Ambiguous
  - Too technical
  - Unrepresentative of image content

Ambiguous class names



Both named 'bow'

Both named 'bat'

Technical class names

Class name:
2007 Cadillac Escalade EXT Crew Cab

Class name:
A340-200

# Prompt context learning

- Learn prompt context word embeddings (frozen vision-language)

Data efficient
- ☑ Improves over zero-shot performance
- ☑ Address prompt sensitivity limitations
- ☑ Maintain open-set properties

- ☒ Relies on handcrafted class names
- ☒ Difficult continual adaptation
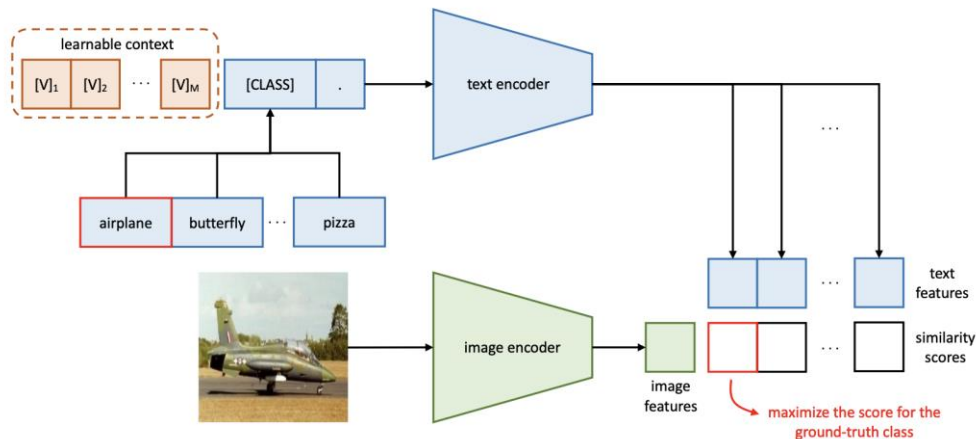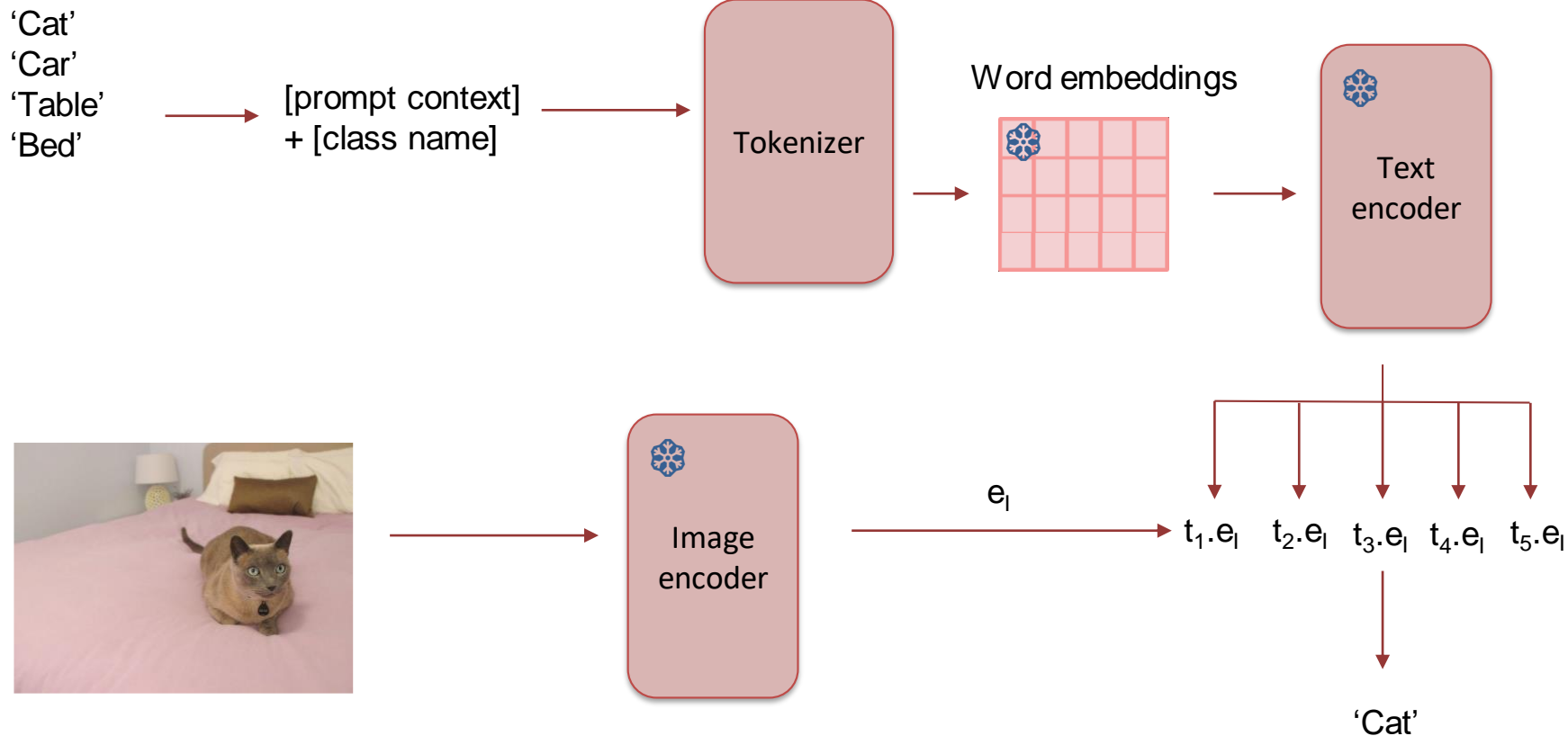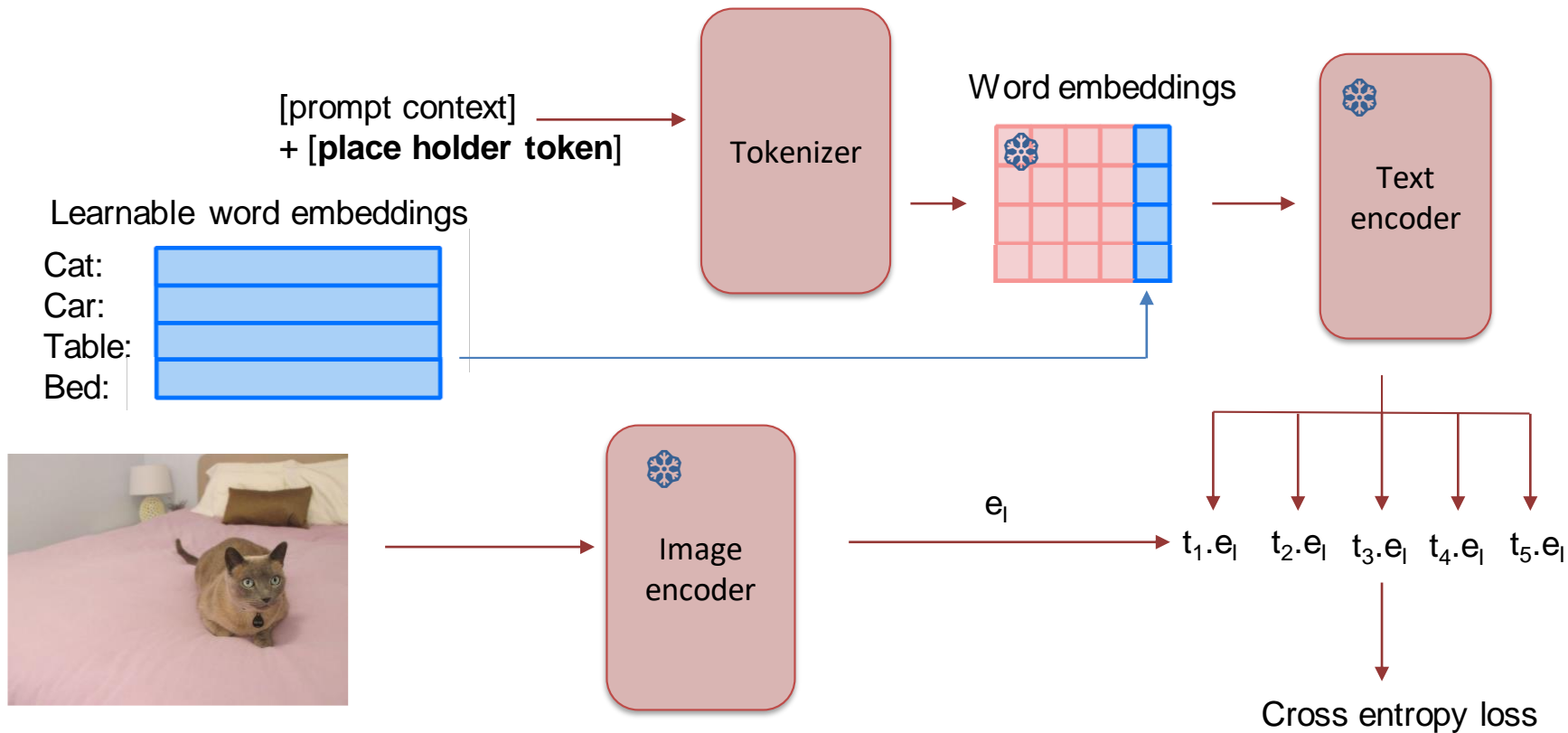- ☒ Weak object detection performance
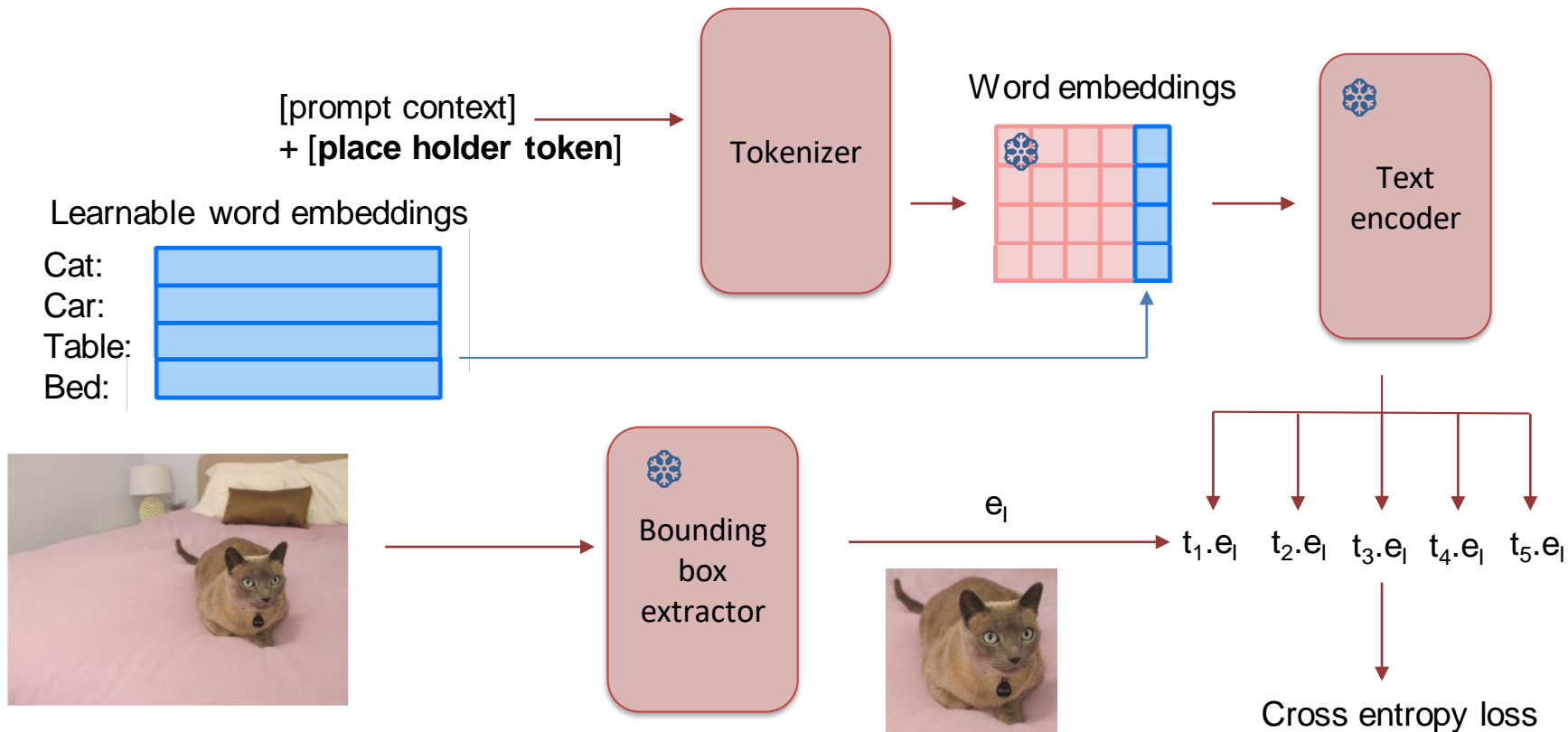


Figure 2: Overview of context optimization (CoOp).

Zhou et al. "Learning to prompt for vision-language models." *International Journal of Computer Vision* (2022)

# Proposed solution

'Cat'
'Car'
'Table'
'Bed'

→ [prompt context] + [class name] → Tokenizer → Word embeddings → Text encoder

$t_1.e_l$  $t_2.e_l$  $t_3.e_l$  $t_4.e_l$  $t_5.e_l$

Image encoder → $e_l$ →

'Cat'

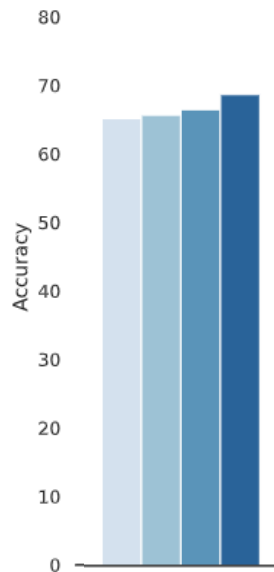# Proposed solution

# Proposed solution

# Experiments: classification with CLIP

- Outperforms SOTA in open-vocabulary and sequential training settings

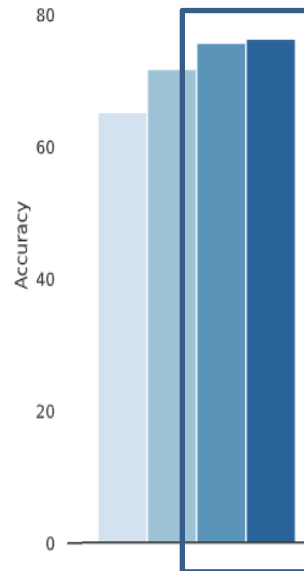- Learning all class names strongly reduces dependency on prompt context

**Open-vocabulary setting**: learning half of the dataset class names

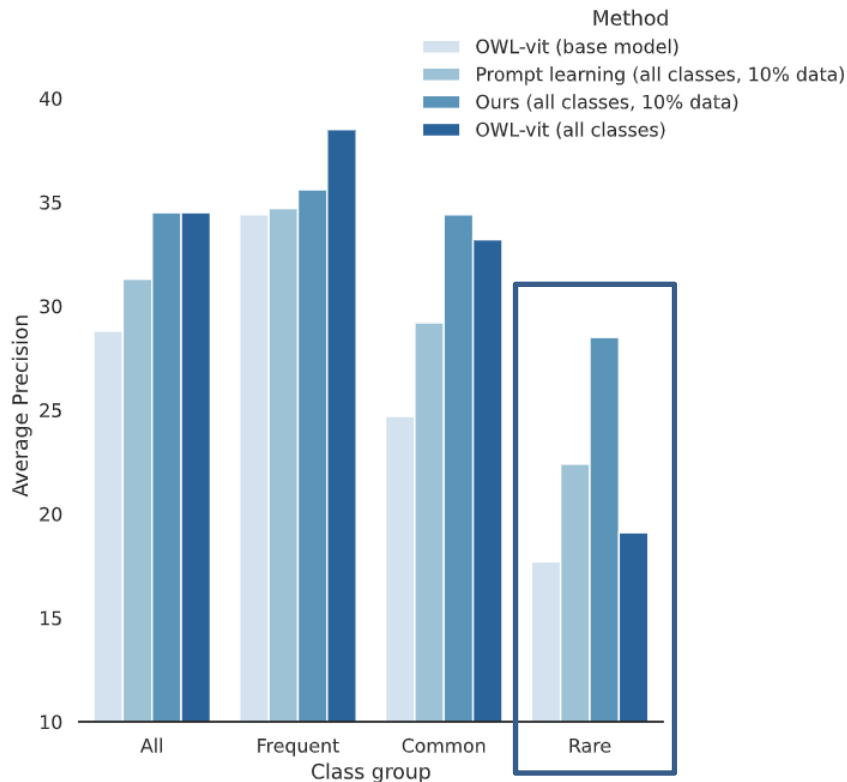**Sequential training setting**: learning two sets of class names sequentially

Method - * with engineered context
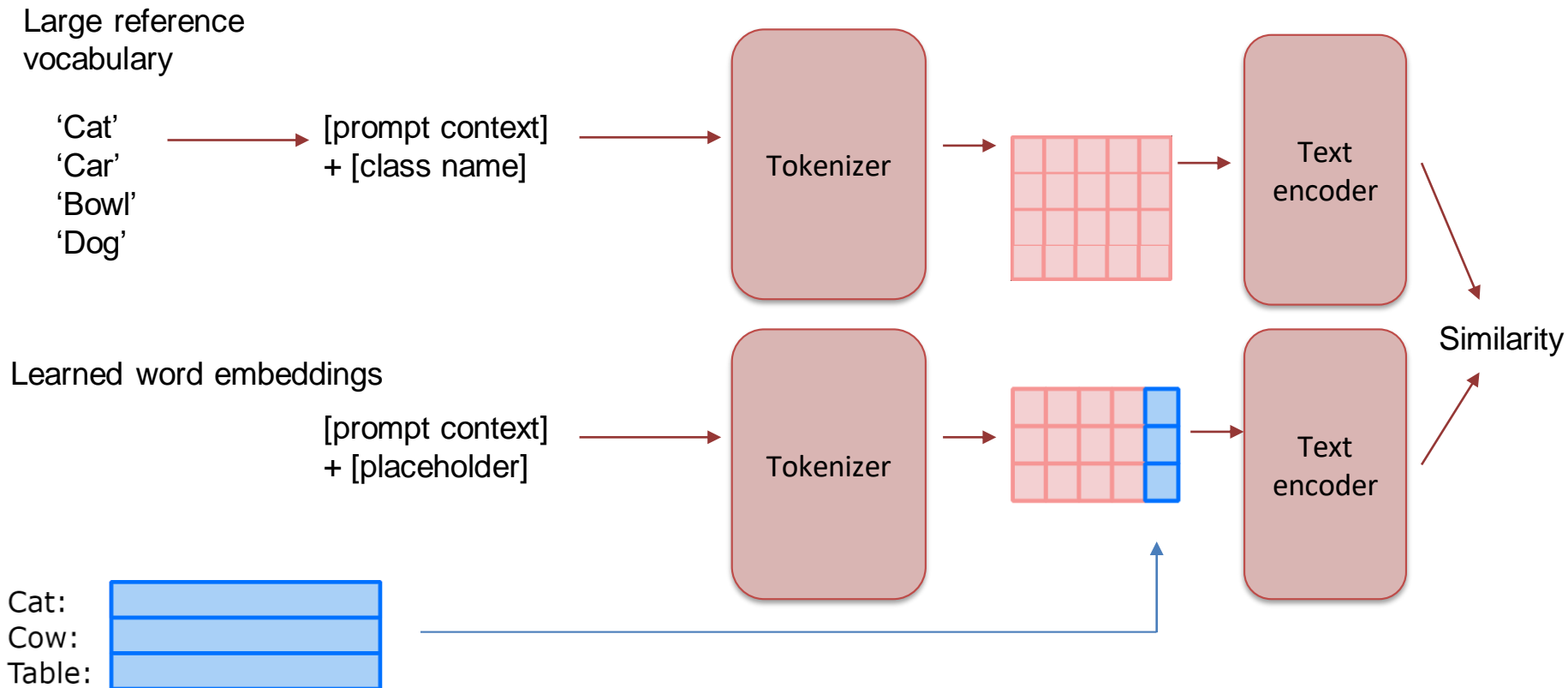
- CLIP*
- CoOp
- Ours
- Ours*

# Experiments: Object detection with OWL-vit

- Learning class names (10% of data) – match performance of fully fine-tuned model

- Significant performance improvement for rare classes

- Significant gains compared to prompt context learning

# Interpretability

Large reference vocabulary

'Cat'
'Car'
'Bowl'
'Dog'

[prompt context] + [class name]

Tokenizer

Text encoder

Learned word embeddings

[prompt context] + [placeholder]

Tokenizer

Text encoder

Similarity

Cat:
Cow:
Table:

# Interpretability



*Original name*:
Arctic

*Original name*:
Tricycle

*Original name*:
Miscellaneous

Boot, ski boot

Cart, rickshaw

wheel,  waterwheel

# Interpretability

Identifying model biases: American English over British English

Original name:

Clothes hamper

↓

Laundry basket

Original name:

Wall socket

↓

Power outlet

Original name:
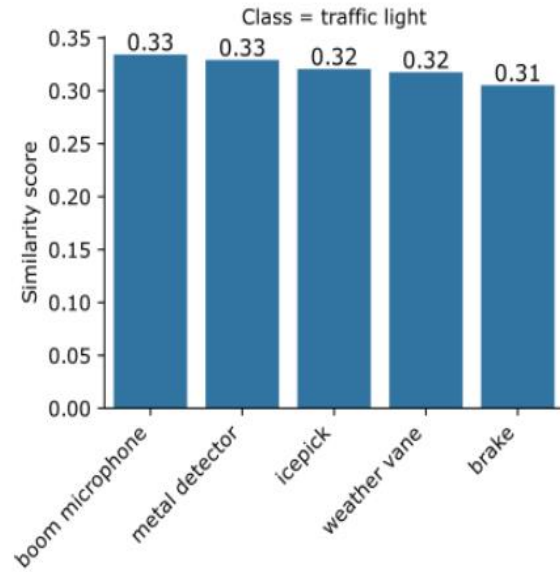
Postbox

↓

Mailbox

Original name:

Trousers

↓

Clothes, Pants

# Interpretability

Potential to identify mislabelled data and failures modes of our method



class examples

# Conclusion

- Novel data efficient adaptation for vision-language models
  - Removes dependency on hand-crafted class names
  - Learn optimal class word embeddings from visual content

- Out of the box usage on classification, detection models

- Complementary to prompt context learning methods

- High interpretability

THU-PM-274

JUNE 18-22, 2023
# CVPR
VANCOUVER, CANADA

# Learning to Name classes for Vision and Language Models

Sarah Parisot, Yongxin Yang, Steven McDonagh

Paper:

Huawei Noah's Ark Lab