# CS 489/698: Introduction to Natural Language Processing

# Lecture 8: Grounded Semantics & Multimodal Language Models

Instructor: Freda Shi

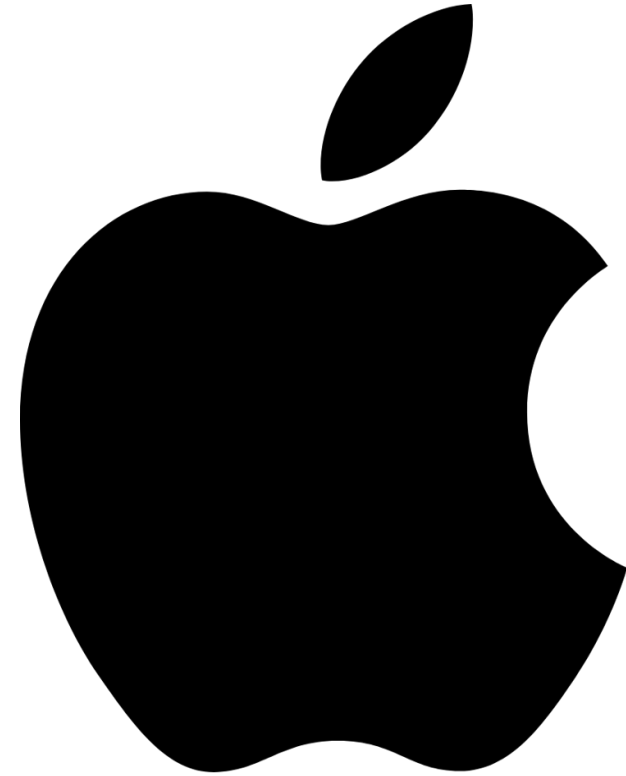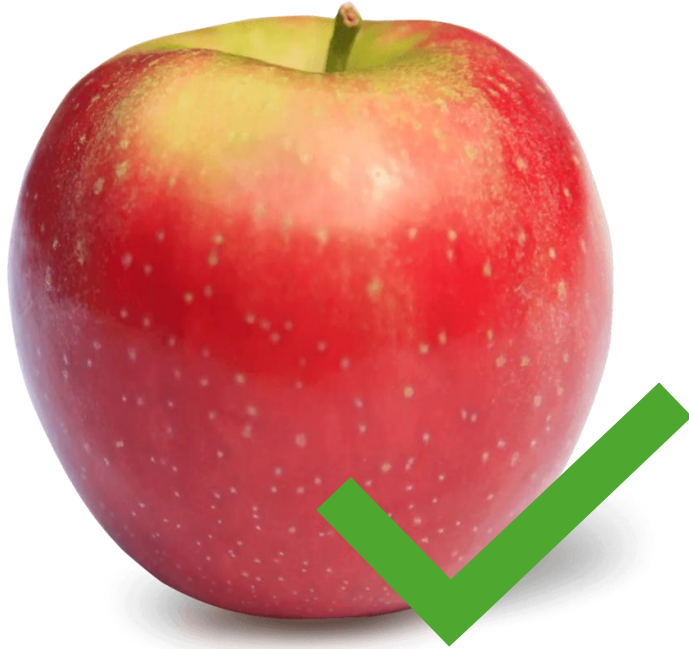*fhs@uwaterloo.ca*

*February 11th, 2026*

UNIVERSITY OF
**WATERLOO**

# Outline

- **Grounded Semantics**
  - **The symbol grounding problem and what *grounding* is**
- Vision-language models (VLMs)
  - Visual-semantic embeddings and CLIP
  - Generative vision-language models
  - Tasks and limitations of VLMs
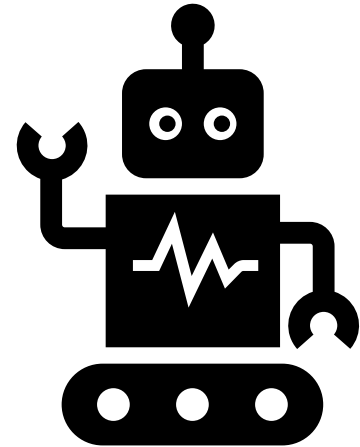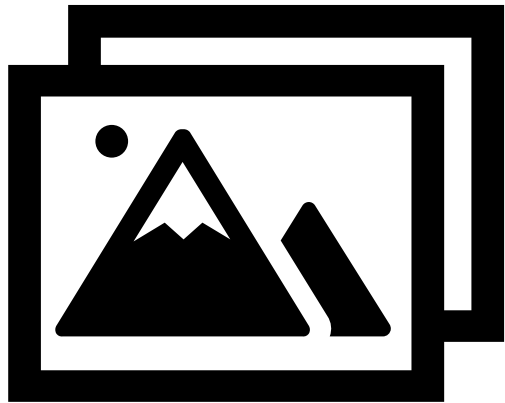
# Meaning in the Real World

My favourite fruit is apple.

# Experience Grounds Language

- Bisk et al. (2020):

  *We posit that the present success of representation learning approaches trained on large, text-only corpora requires the parallel tradition of research on the broader physical and social context of language to address the deeper questions of communication.*

[Bisk et al. 2020. *Experience Grounds Language. In*: EMNLP]

# Grounded Semantics

- Meanings demonstrated from other sources of data in addition to the language systems.

## Distributional Semantics

A bottle of *tezgüino* is on the table.
Everybody likes *tezgüino*.
Don't have *tezgüino* before you drive.
We make *tezgüino* out of corn.

## Visually Grounded Semantics

*tezgüino* = 

[Figure credit: Alejandro Linares Garcia]

# The Symbol Grounding Problem (Harnad, 1990)

- Symbol → meaning: how to make sense of symbols?

- Practical implication: enable the reliably *meaningful* interaction between language models and humans/physical world.
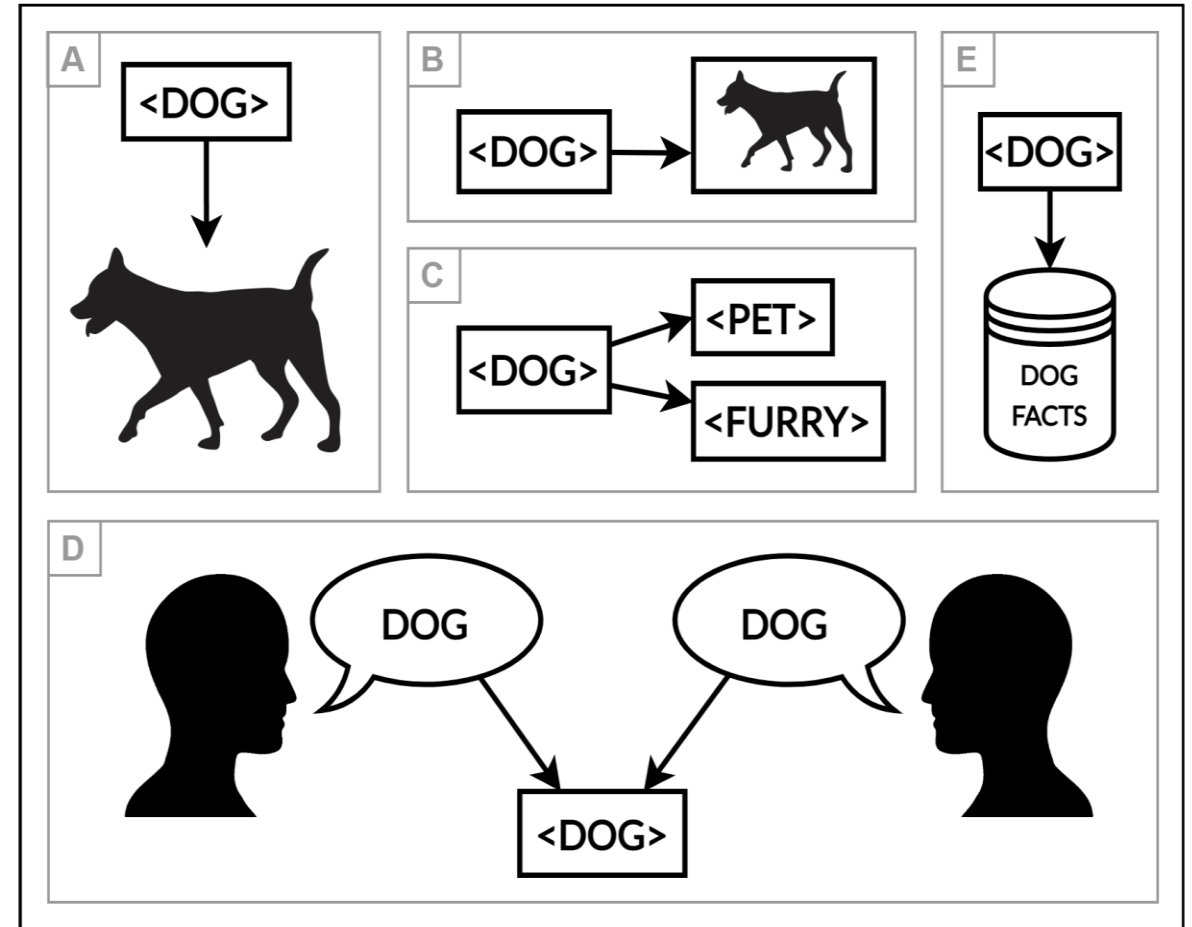


Stochastic parrots or semantic comprehension?

Still under debate...

But we all agree -- external source of meaning (e.g., data from another modality) better implies comprehension.

[Bender et al. 2021. *On the Dangers of Stochastic Parrots: Can Language Models Be Too Big? In: ACM FAccT*; Figure source: https://freesvg.org]

# Taxonomy of Grounding

- Grounding can be categorized into
  - A: referential grounding
  - B: sensorimotor grounding
  - C: relational grounding
  - D: communicative grounding
  - E: epistemic grounding

- Chai et al. (2018)
  - A, B, C, E: semantic (static) grounding
  - D: communicative (dynamic) grounding



[Mollo and Millière. 2023. The Vector Grounding Problem]

# Recap: Text-Only Language Models

- Two popular types of text-only (ungrounded) language models:
  - Autoregressive models (e.g., GPT; Radford et al., 2018) – better for generation

$$P_\Theta(w_{n+1} \mid w_1, \ldots, w_n)$$

  - Masked language models (e.g., BERT; Devlin et al., 2019) – better for feature extraction

$$P_\Theta(w_i \mid w_1, \ldots, w_{i-1}, w_{i+1}, w_n)$$

- Incorporating visual signals leads to two families of vision-language models:
  - BERT → Joint visual semantic embeddings.
  - GPT → Generative vision-language models.
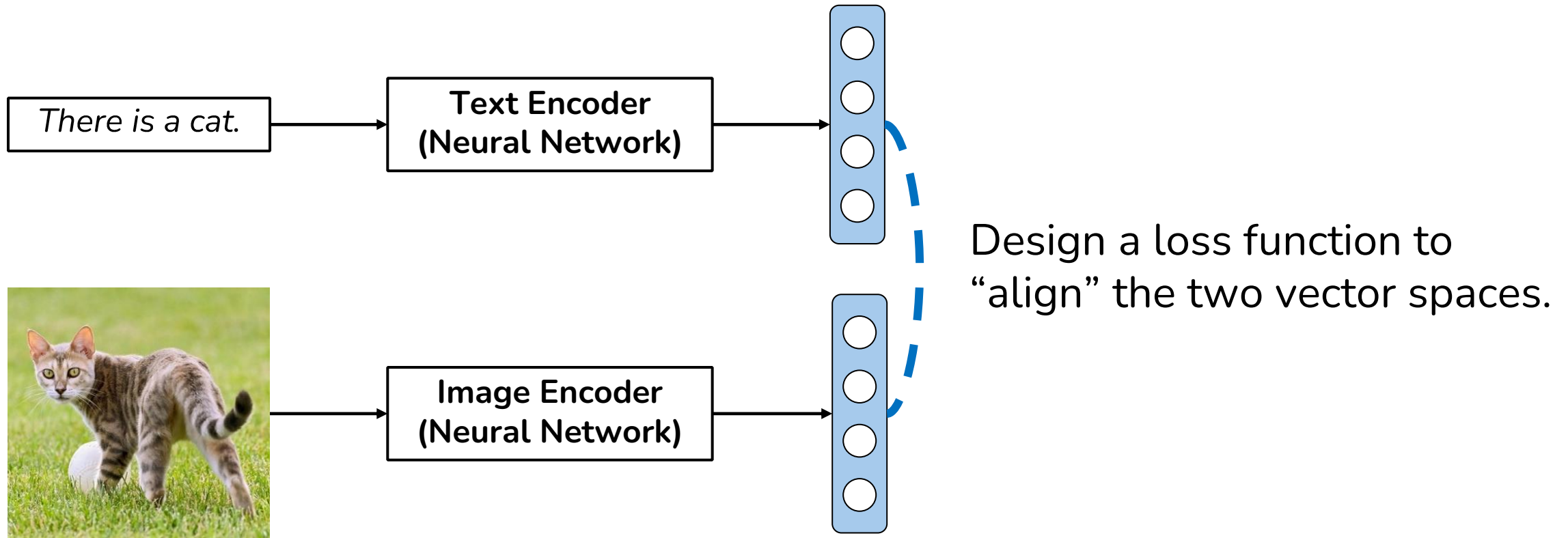
# Outline

- Grounded Semantics

  - The symbol grounding problem and what *grounding* is

- **Vision-language models (VLMs)**

  - **Visual-semantic embeddings and CLIP**

  - Generative vision-language models

  - Tasks and limitations of VLMs

# Joint Visual-Semantic Embedding Space

Idea: encode visual and textual information into a shared space.

Embedding: vector space.



Design a loss function to "align" the two vector spaces.

# Joint Visual-Semantic Embedding Space: Dataset

Training data: images and their text descriptions.

Example: Microsoft COCO (Lin et al., 2014) collects 80K images of common objects and their captions.



a street sign on the corner of a busy street
a city street scene with a street sign, "church st." in view.
a city street filled with lots of traffic surrounded by tall buildings.
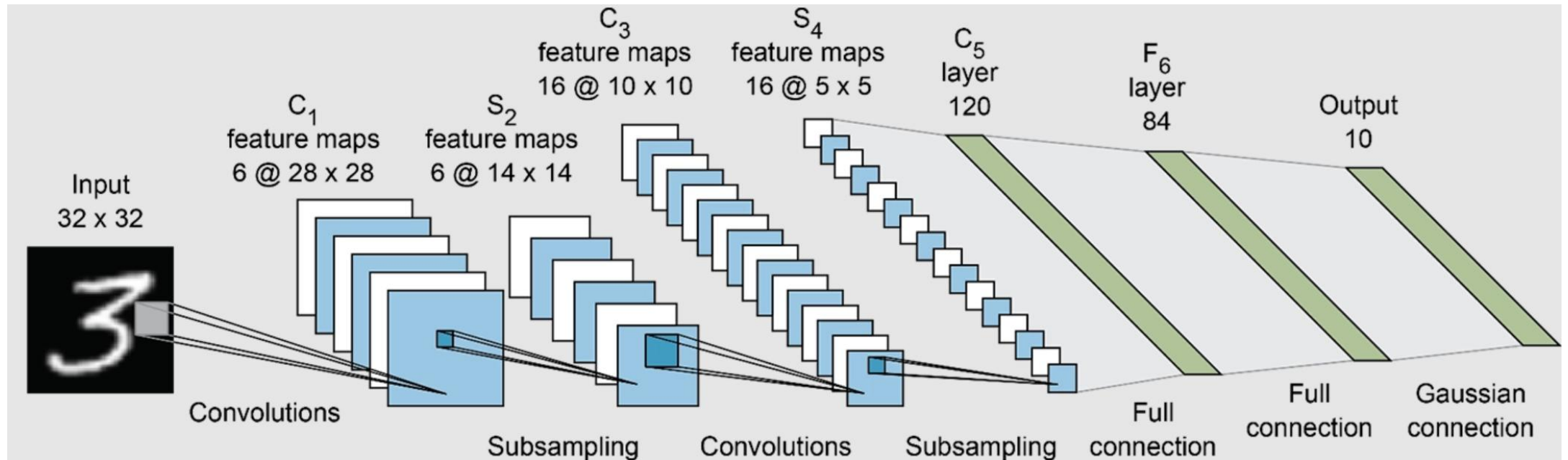busy city street with cars parked on the side of the road.
a street sign says church street in a city

# Visual Encoders

Convert an image to a fixed-dimensional vector representation.

# Joint Visual-Semantic Embedding Space: Objective

Core idea: Matched image-caption pair should be closer than mismatched pairs in the embedding space.



There is a cat.     There is an apple.

model parameter     "margin"

$$\Theta^* = \text{argmin}_\Theta \sum_{(I^+, T^+, T^-)} \max\left(0, \delta - \cos(I_\Theta^+, T_\Theta^+) + \cos(I_\Theta^+, T_\Theta^-)\right)$$
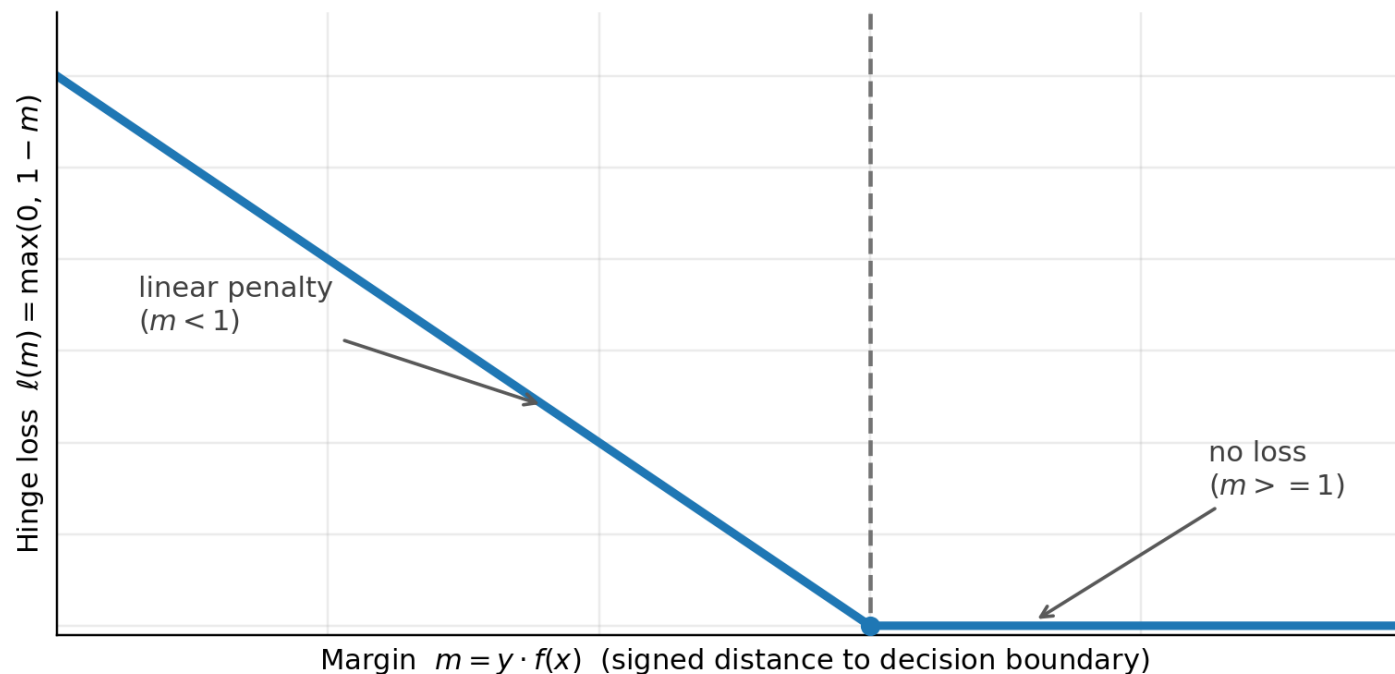
enumerate over dataset

$$+ \sum_{(T^+, I^+, I^-)} \max\left(0, \delta - \cos(T_\Theta^+, I_\Theta^+) + \cos(T_\Theta^+, I_\Theta^-)\right)$$

"Triplet-Based Hinge Loss"

[Kiros et al. 2014. Unifying visual-semantic embeddings with multimodal neural language models.]

# "Hinge" Loss

$$\Theta^* = \text{argmin}_\Theta \sum_{(I^+, T^+, T^-)} \max\left(0, \delta - \cos(I_\Theta^+, T_\Theta^+) + \cos(I_\Theta^+, T_\Theta^-)\right)$$

$$+ \sum_{(T^+, I^+, I^-)} \max\left(0, \delta - \cos(T_\Theta^+, I_\Theta^+) + \cos(T_\Theta^+, I_\Theta^-)\right)$$



linear penalty
($m < 1$)

no loss
($m >= 1$)

Hinge loss $\ell(m) = \max(0, 1 - m)$

Margin $m = y \cdot f(x)$ (signed distance to decision boundary)

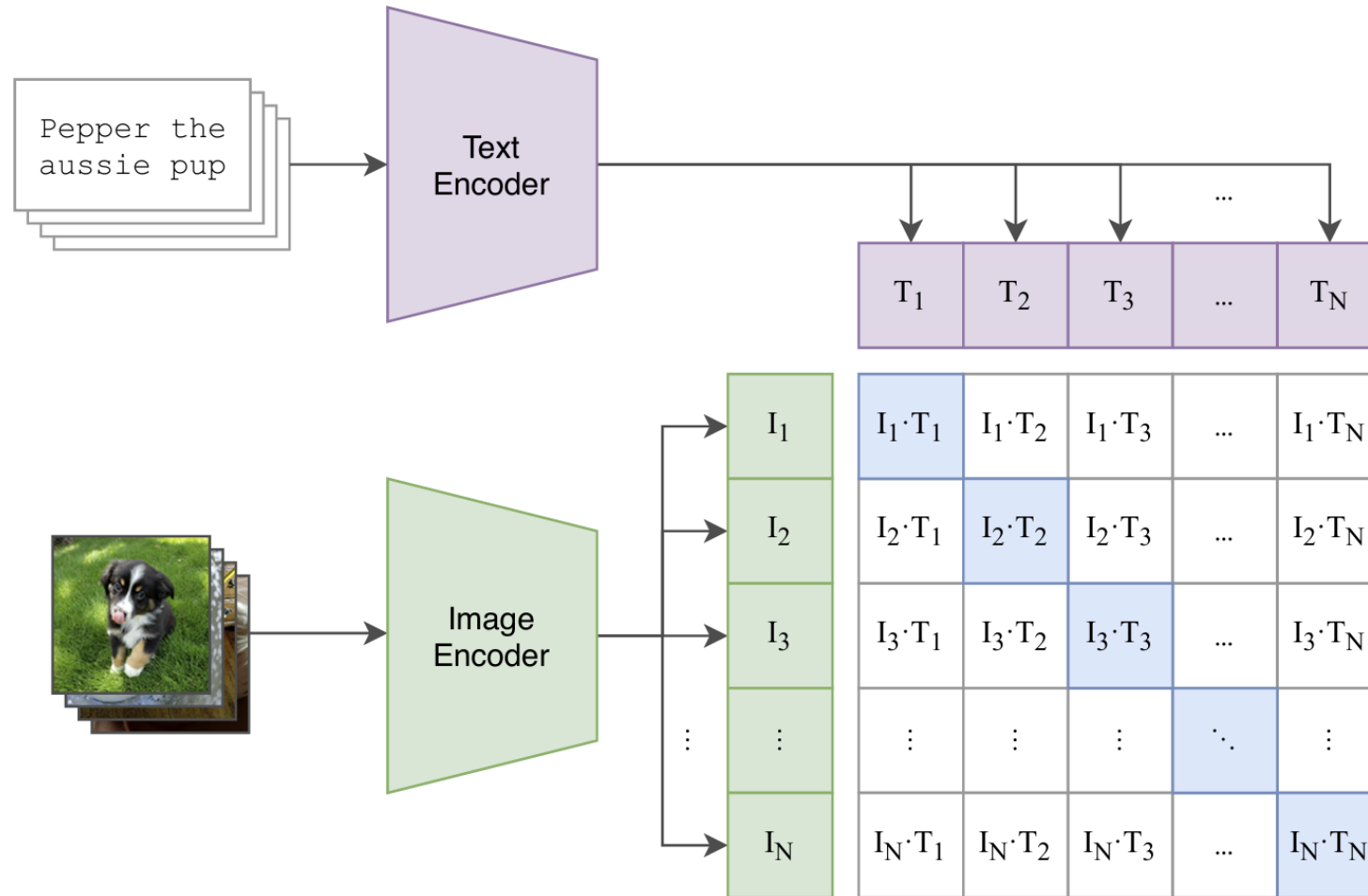# Properties of the Joint Space

Images and text are close in a good joint embedding space if they are semantically related.

**Example applications:**

- Bidirectional image-caption retrieval, e.g., Google image search.

- Image captioning

Text in the training data can be at any level of granularity (words, phrases, sentences, paragraph, documents, etc.).

# Contrastive Language-Image Pretraining (CLIP)



[Figure: Radford et al. 2021. Learning transferable visual models from natural language supervision]

# CLIP Objective: "Classifying for the True Label"

- Image-to-text retrieval: given a pool of text, model the probability of choosing the correct text; and vice versa.

$$\Theta^* = \arg\min_{\Theta} \mathbb{E}_{[(I_1,T_1),\ldots(I_n,T_n)]} \left[ \sum_i -\log P_\Theta(T_i \mid I_i; [T_{1\ldots n}]) - \log P_\Theta(I_i \mid T_i; [I_{1\ldots n}]) \right]$$

text pool
$n \times d$-dimensional features

image pool
$n \times d$-dimensional features



$$P_\Theta\left(T_j \mid I_i; [T_{1..n}]\right) = \mathsf{softmax}\left([T_{1..n}] \cdot I_i\right)_j = \frac{\exp(\langle I_i, T_j\rangle)}{\sum_{k=1}^n \exp(\langle I_i, T_k\rangle)}$$
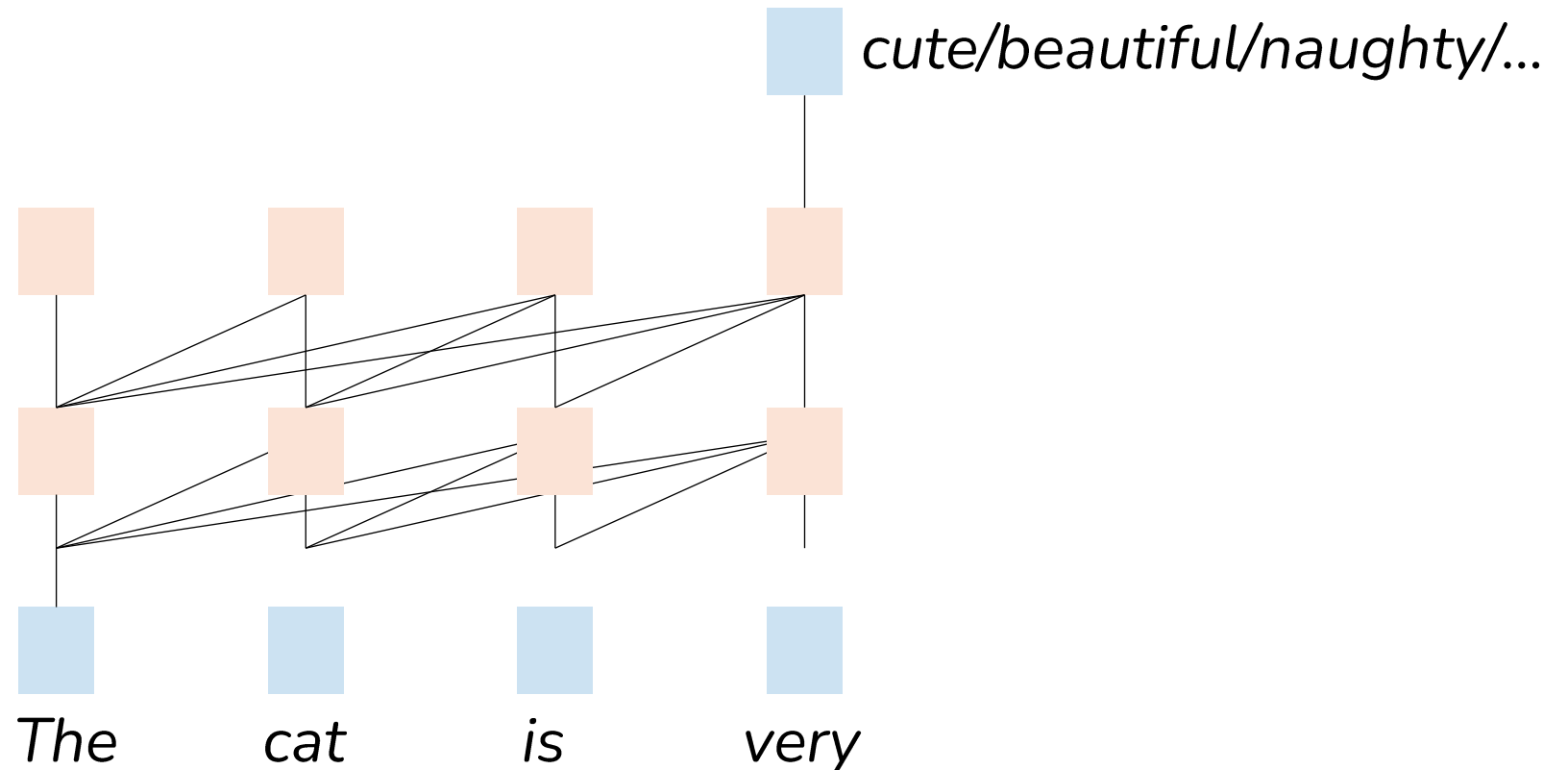
There is a cat. $(T_i)$

# Outline

- Grounded Semantics
  - The symbol grounding problem and what *grounding* is
- **Vision-language models**
  - Visual-semantic embeddings and CLIP
  - **Generative vision-language models**
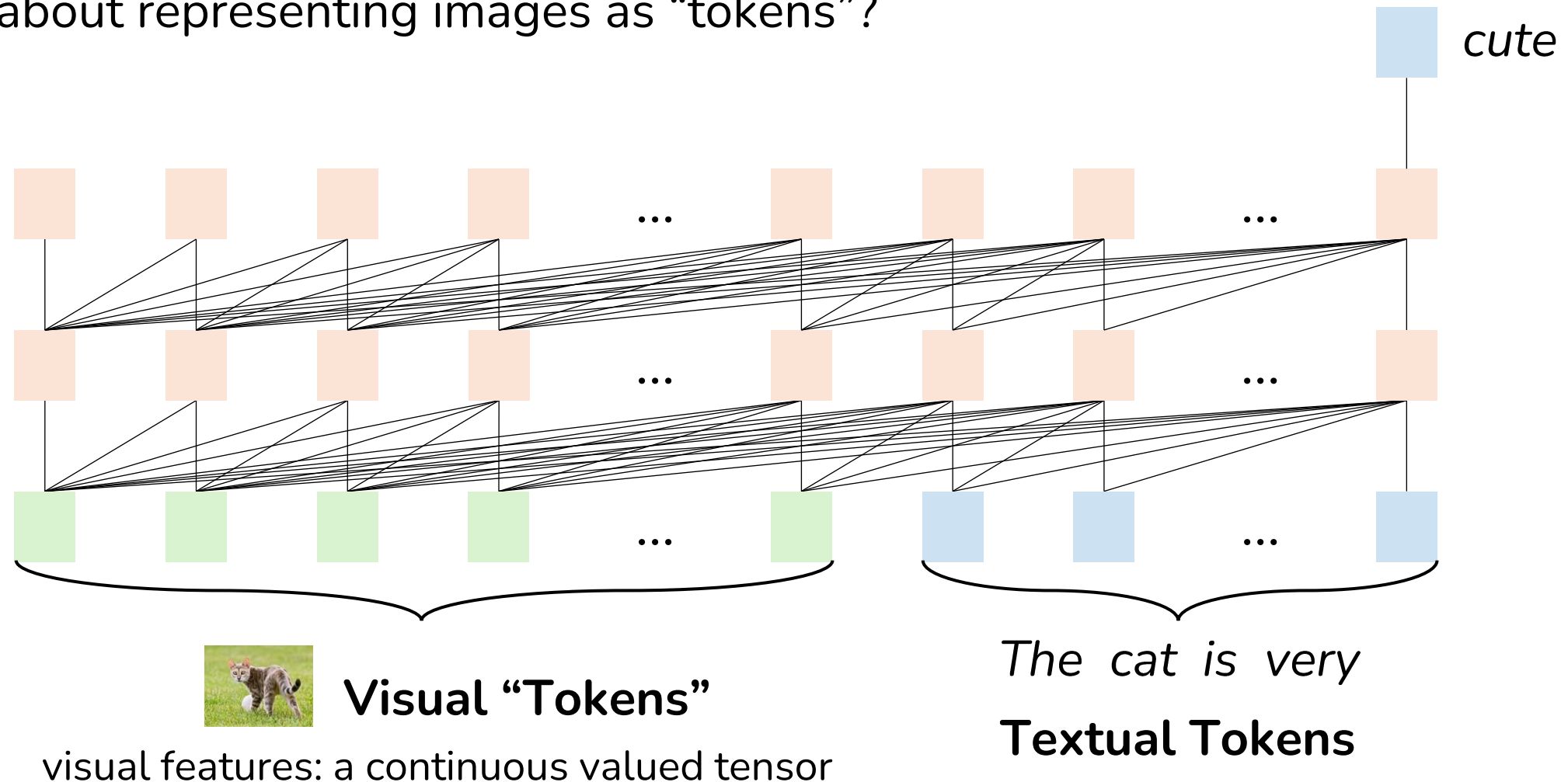  - Tasks and limitations of VLMs

# Recap: Generative Autoregressive Language Models

Text-only language models: predicting the next token conditioned on the history.



*cute/beautiful/naughty/...*

*The*　　*cat*　　*is*　　*very*

# Extending to Vision-Language Models (VLMs)

How about representing images as "tokens"?



*cute*

*The cat is very*

**Visual "Tokens"**

visual features: a continuous valued tensor

**Textual Tokens**

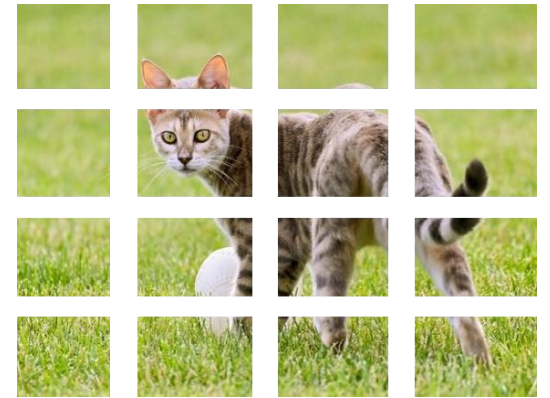[Liu et al. 2023. *Visual instruction tuning. In:* NeurIPS.]

# Generative VLM: Training Objective

$$\Theta^* = \min_{\Theta} \sum_i \sum_j -\log P_{\Theta}(w_j^{(i)}; \mathsf{Image}^{(i)}, w_j^{(i)}, \ldots, w_{j-1}^{(i)})$$

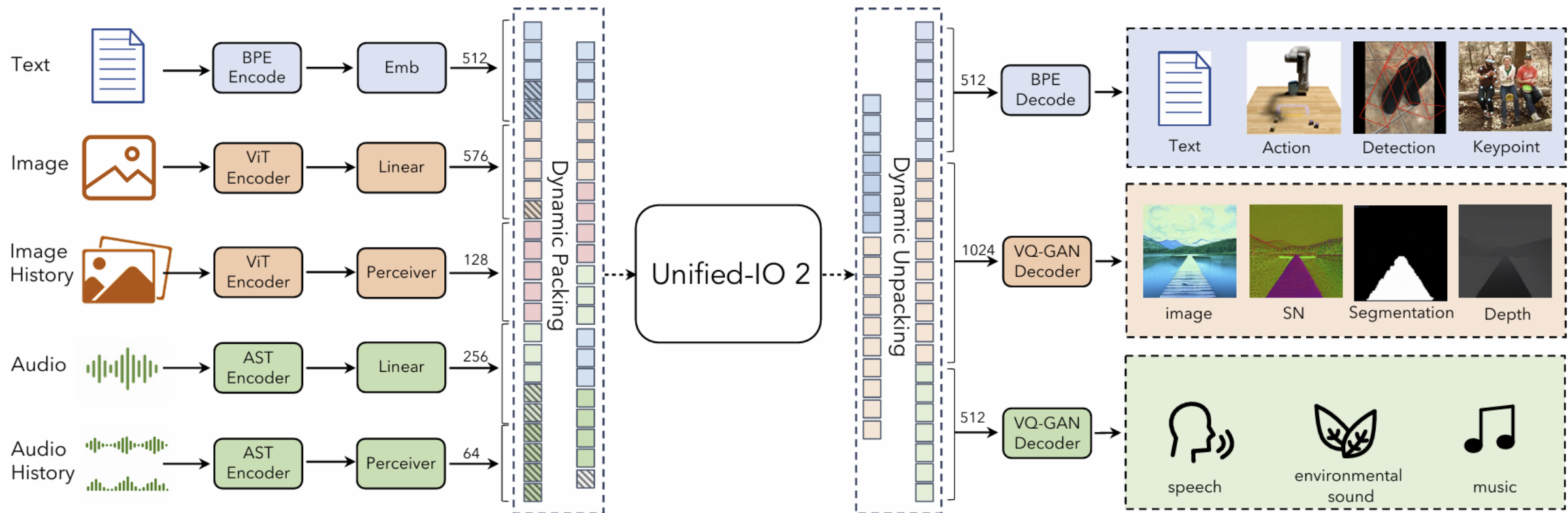Loss function only calculated on textual positions.



In practice, each visual token correspond to an image patch.

Visual encoders use patches to improve representation quality.

Training all involved parameters via backpropagation and gradient descent.
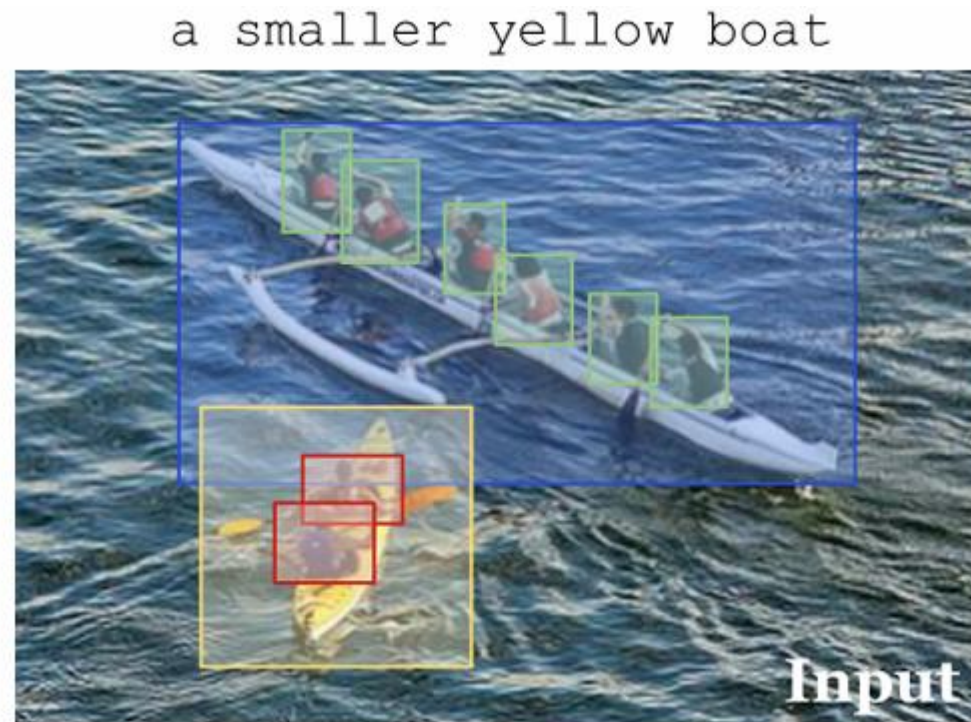
# Towards Encoding Everything in the World



[Lu et al. 2024. *Unified-IO 2: Scaling autoregressive multimodal models with vision, language audio and action. In:* CVPR.]

# Outline

- Grounded Semantics

  - The symbol grounding problem and what *grounding* is

- **Vision-language models**

  - Visual-semantic embeddings and CLIP

  - **Generative vision-language models**

  - Tasks and limitations of VLMs

# Finer-Grained Vision-Language Tasks

- **Object retrieval** (assuming all objects' bounding boxes are given).
  - Cognitive plausibility: recognizing objects are very easy for humans (in fact, 5-month-old infants; Baillargeon et al., 1985).

# Finer-Grained Vision-Language Tasks

- **Multimodal coreference resolution** (w/o assuming bounding boxes)
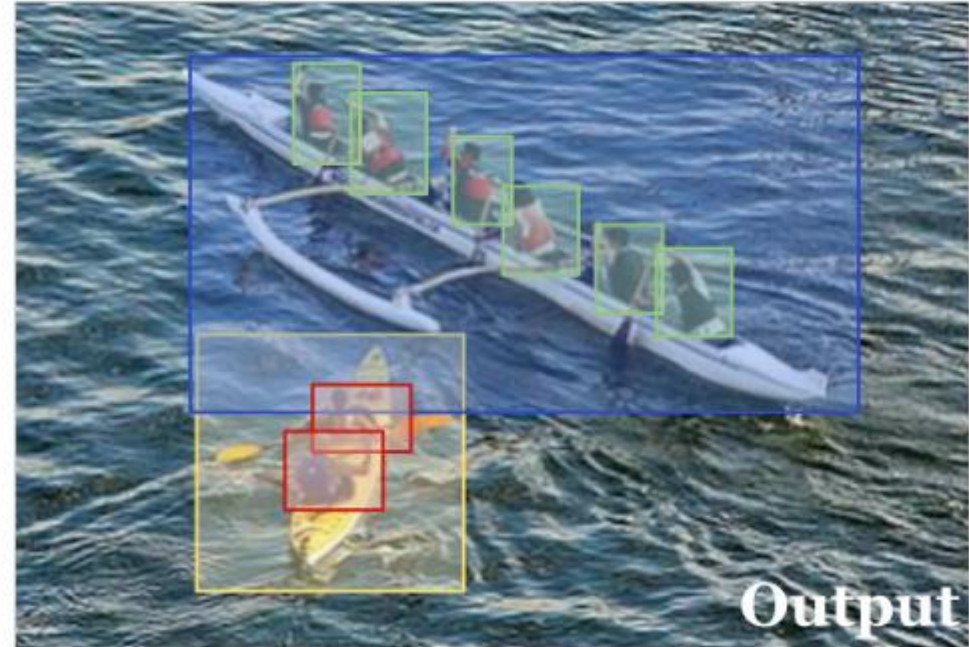
# Finer-Grained Vision-Language Tasks

- **Phrase grounding:** mapping phrases to objects in the image.

- **Dense captioning (reverse):** write a short description for each detected object.

# Limitations of Current VLMs

- Lack of physical knowledge, and the neural architecture makes it hard to incorporate the knowledge.



[Sarkar et al. 2024. *Shadows don't lie and lines can't bend! Generative models don't know projective geometry... for now. In:* CVPR.]

# Limitations of Current VLMs

- Poor in recognizing spatial relations.



*The tree is behind the car.*

*The tree is to the right of the car.*

*The tree is in front of the car.*

*The car is to the left of the tree.*

[Zhang et al. 2025. *Do vision-language models represent space and how? Evaluating spatial frame of reference under ambiguities. In:* ICLR]
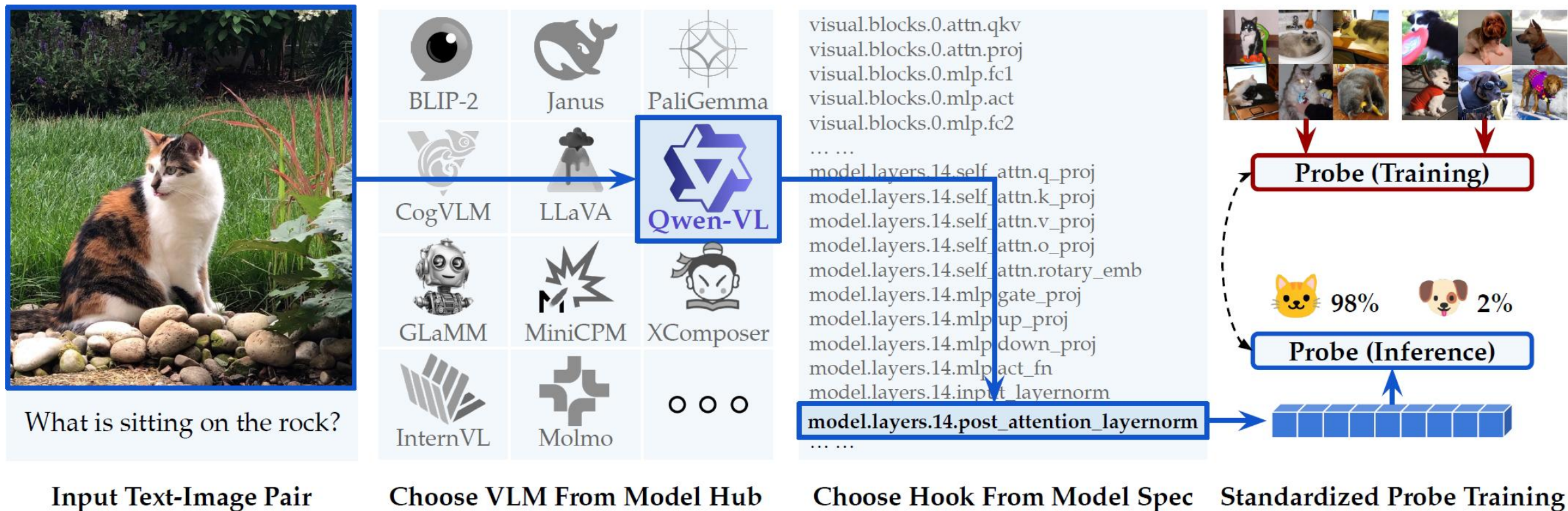
# Limitations of Current VLMs

Lack of cultural diversity representation.



[Bhatia et al. 2024. *From Local Concepts to Universals: Evaluating the Multicultural Understanding of Vision-Language Models. In:* EMNLP]

# Analyzing Internals of VLMs



**Input Text-Image Pair** · **Choose VLM From Model Hub** · **Choose Hook From Model Spec** · **Standardized Probe Training**

https://github.com/compling-wat/vlm-lens

[Sheta et al. 2025. *From Behavioral Performance to Internal Competence: Interpreting Vision-Language Models with VLM-Lens. In:* EMNLP Systems Demonstraton]

# Next

- Assignment 2 will be released on Friday

- Victor takes over the lectures on pretraining language models