# Grounding in Images

‣ How would you describe this image?

‣ What does the word "*spoon*" evoke?



*the girl is licking the spoon of batter*

# Grounding Spoon



Winco 0005-03 7 3/8" Dinner Spoon...

$7.16



wikiHow

How to Hold a Spoon: 13 Steps (...



GO Indiegogo

Spoon that Elevates Taste ...

# Grounding Language in Images

‣ Syntactic categories have some regular correspondences to the world:

   ‣ Nouns: objects

   ‣ Verbs: actions

   ‣ Sentences: whole scenes or things happening

‣ Tasks:

   ‣ Object recognition (pick out one most salient object or detect all of them)

   ‣ Image captioning: produce a whole sentence for an image

# Language-vision Models
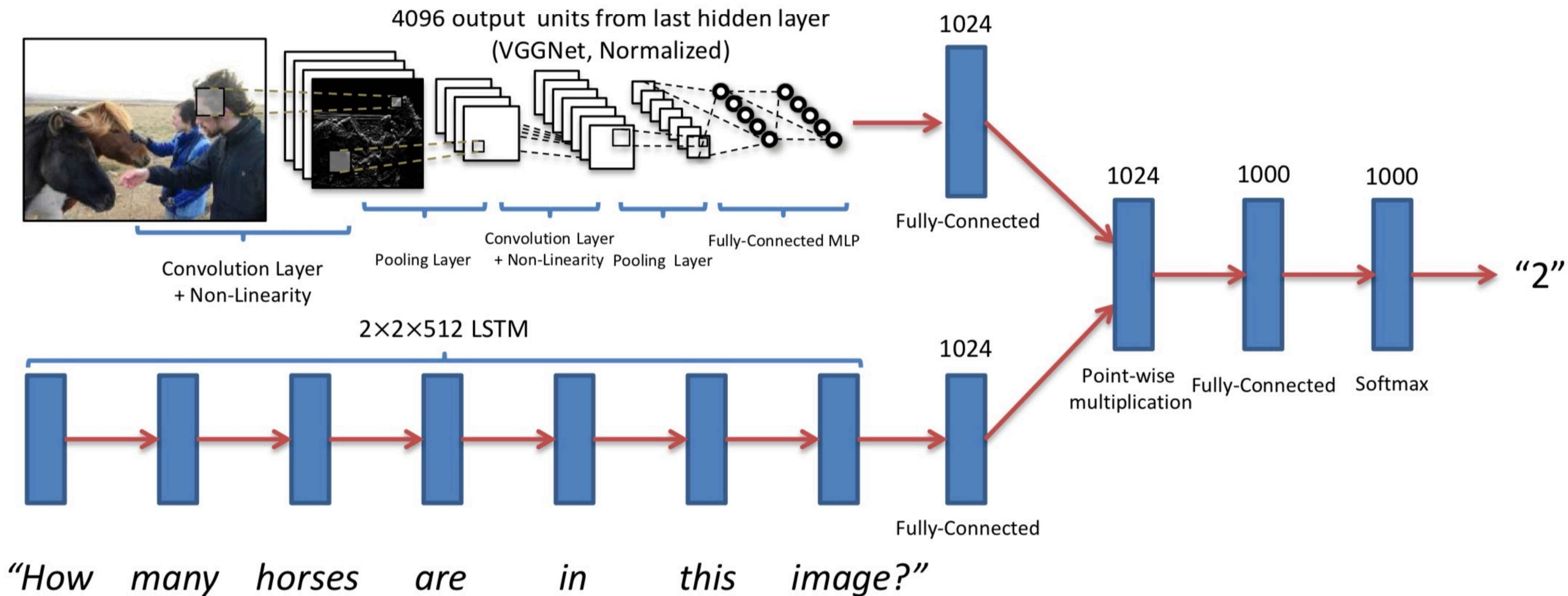


Image encoder
(CNN, Transformer)

*the girl is licking the spoon of batter*

Language encoder
(LSTM, Transformer)
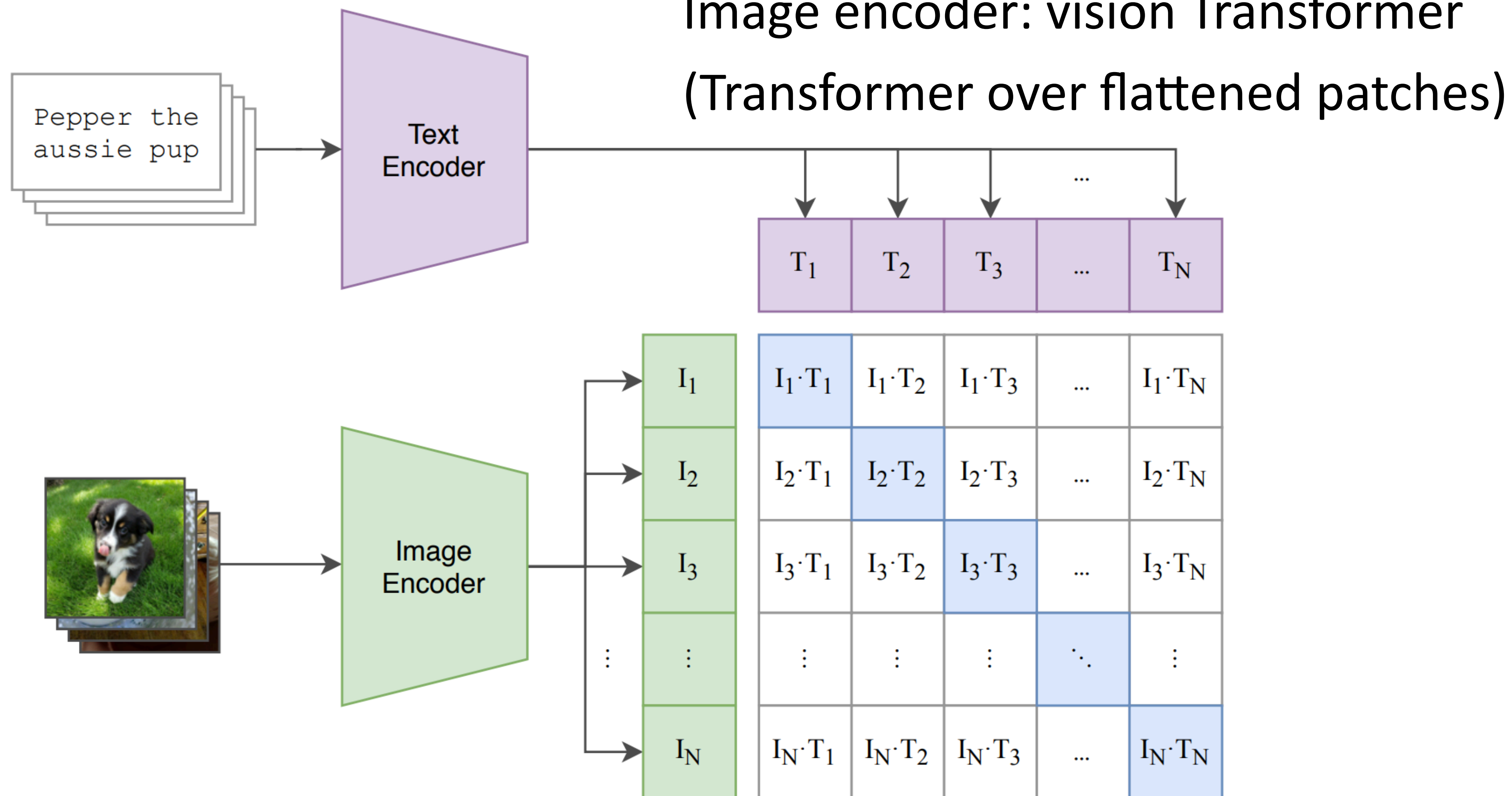
Cross-attention/joint layer

Prediction

# Visual Question Answering



Agrawal et al., 2015

# Language-vision Pre-training: CLIP

(1) Contrastive pre-training

Text encoder: Transformer

Image encoder: vision Transformer

(Transformer over flattened patches)



Radford et al., 2021

# Language-vision Pre-training: CLIP

| | $T_1$ | $T_2$ | $T_3$ | ... | $T_N$ |
|---|---|---|---|---|---|
| $I_1$ | $I_1 \cdot T_1$ | $I_1 \cdot T_2$ | $I_1 \cdot T_3$ | ... | $I_1 \cdot T_N$ |
| $I_2$ | $I_2 \cdot T_1$ | $I_2 \cdot T_2$ | $I_2 \cdot T_3$ | ... | $I_2 \cdot T_N$ |
| $I_3$ | $I_3 \cdot T_1$ | $I_3 \cdot T_2$ | $I_3 \cdot T_3$ | ... | $I_3 \cdot T_N$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\ddots$ | $\vdots$ |
| $I_N$ | $I_N \cdot T_1$ | $I_N \cdot T_2$ | $I_N \cdot T_3$ | ... | $I_N \cdot T_N$ |

‣ Contrastive objective: each image should be more similar to its correspond caption than to other captions

$$\text{maximize softmax}(I_1^T T_i)[1]$$
$$+ \text{softmax}(I_2^T T_i)[2]$$
$$+ \dots$$

Radford et al., 2021

# Language-vision Pre-training: CLIP

(2) Create dataset classifier from label text

| plane |
|-------|
| car |
| dog |
| ⋮ ⋮ |
| bird |

A photo of a {object}.

Text Encoder

| $T_1$ | $T_2$ | $T_3$ | ... | $T_N$ |
|-------|-------|-------|-----|-------|

(3) Use for zero-shot prediction

Image Encoder

$I_1$

| $I_1 \cdot T_1$ | $I_1 \cdot T_2$ | $I_1 \cdot T_3$ | ... | $I_1 \cdot T_N$ |
|-----------------|-----------------|-----------------|-----|-----------------|

A photo of a dog.

Radford et al., 2021

# CLIP: Zero-shot Results



**Stanford Cars**

correct label: 2012 Honda Accord Coupe    correct rank: 1/196    correct probability: 63.30%

# CLIP: Zero-shot Results



**Country211**

correct label: Belize    correct rank: 5/211    correct probability: 3.92%