

# Computer Vision

Multi-modal learning: Captioning and Speaking

---

Tae-Hyun Oh (오태현)  
전자전기공학과  
POSTECH

Slide by Juyong Lee (이주용)

TAs: {Dongmin Choi , Jongha Kim, Juyong Lee, Sungbin Kim} (in alphabetic order)

# 1. Overview of multi-modal learning

## 2. Multi-modal tasks (1) – Visual data & Text

- 2.1 Text embedding
- 2.2 Joint embedding
- 2.3 Cross modal translation
- 2.4 Cross modal reasoning

## 3. Multi-modal tasks (2) – Visual data & Audio

- 3.1 Sound representation
- 3.2 Joint embedding
- 3.3 Cross modal translation
- 3.4 Cross modal reasoning

1.

# Overview of multi-modal learning

Toward the world beyond images

# 1. Overview of multi-modal learning

Multi-modal learning overview

## Modalities in multi-modal learning



Depth



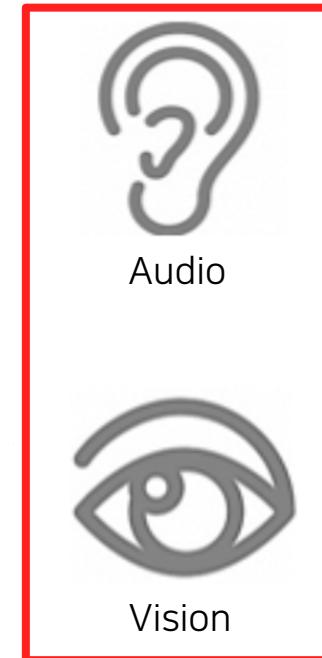
Odor



Unimodal



Texture



Force

. 오늘 8시 미팅은 아래 주소로 하겠습니다.  
교수님, 오늘 전체 조교들이 조교 미팅 들어가  
혹시라도 예정시간 보다 늦어지면 미리 말씀!  
하겠습니다.

PAMI Lab is inviting you to a scheduled Z meeting.

Topic: PAMI Lab's Zoom Meeting  
Time: Jan 11, 2021 08:00 PM Seoul  
Join Zoom Meeting

Text

Unimodal

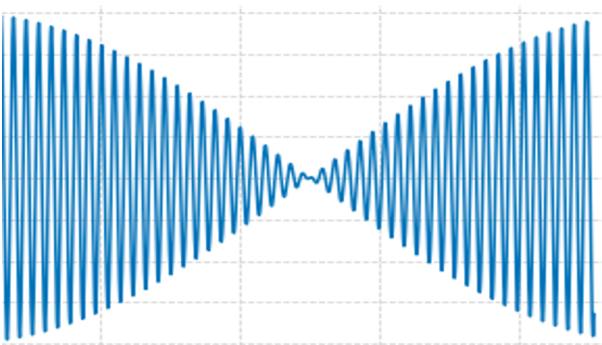
# 1. Overview of multi-modal learning

Multi-modal learning overview

## Challenge (1) - Different representations between modalities



Audio



Image

157	153	174	168	150	152	129	151	172	161	155	156
155	182	163	74	75	62	83	17	110	210	180	154
180	180	50	14	34	6	10	33	48	106	159	181
206	109	5	124	131	111	120	204	166	15	56	180
194	68	197	251	297	299	299	228	227	87	71	201
172	105	207	233	233	214	220	239	228	98	74	206
188	68	179	209	185	215	211	158	139	75	20	169
189	97	165	84	10	168	134	11	31	62	22	148
199	168	191	193	158	227	178	143	182	105	36	190
205	174	155	252	236	231	149	178	228	43	95	234
190	216	116	149	236	187	85	150	79	38	218	241
190	224	147	108	227	210	127	102	36	101	255	224
190	214	173	66	103	143	95	50	2	109	249	215
187	196	236	75	1	81	47	0	6	217	255	211
183	202	237	145	0	0	12	108	200	138	243	236
195	206	123	297	177	121	123	200	175	13	96	218

오늘 8시 미팅은 아래 주소로 하겠습니다.  
교수님, 오늘 전체 조교들이 조교 미팅 들어가  
혹시라도 예정시간 보다 늦어지면 미리 말씀!  
겠습니다.

PAMI Lab is inviting you to a scheduled Z  
meeting.

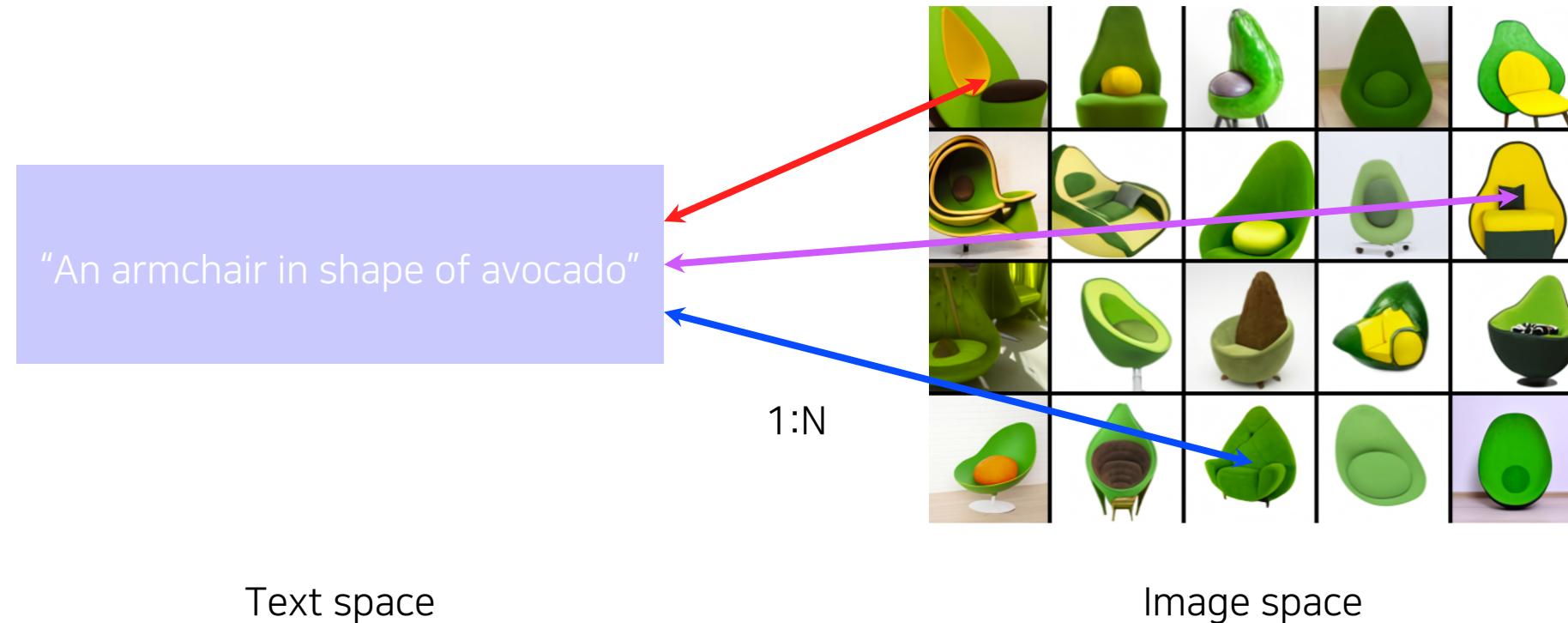
Topic: PAMI Lab's Zoom Meeting  
Time: Jan 11, 2021 08:00 PM Seoul  
Join Zoom Meeting

	living	being	feline	human	gender	royalty	verb	plural
<i>cat</i> →	0.6	0.9	0.1	0.4	-0.7	-0.3	-0.2	
<i>kitten</i> →	0.5	0.8	-0.1	0.2	-0.6	-0.5	-0.1	
<i>dog</i> →	0.7	-0.1	0.4	0.3	-0.4	-0.1	-0.3	
<i>houses</i> →	-0.8	-0.4	-0.5	0.1	-0.9	0.3	0.8	

# 1. Overview of multi-modal learning

Multi-modal learning overview

## Challenge (2) - Unbalance between heterogeneous feature spaces

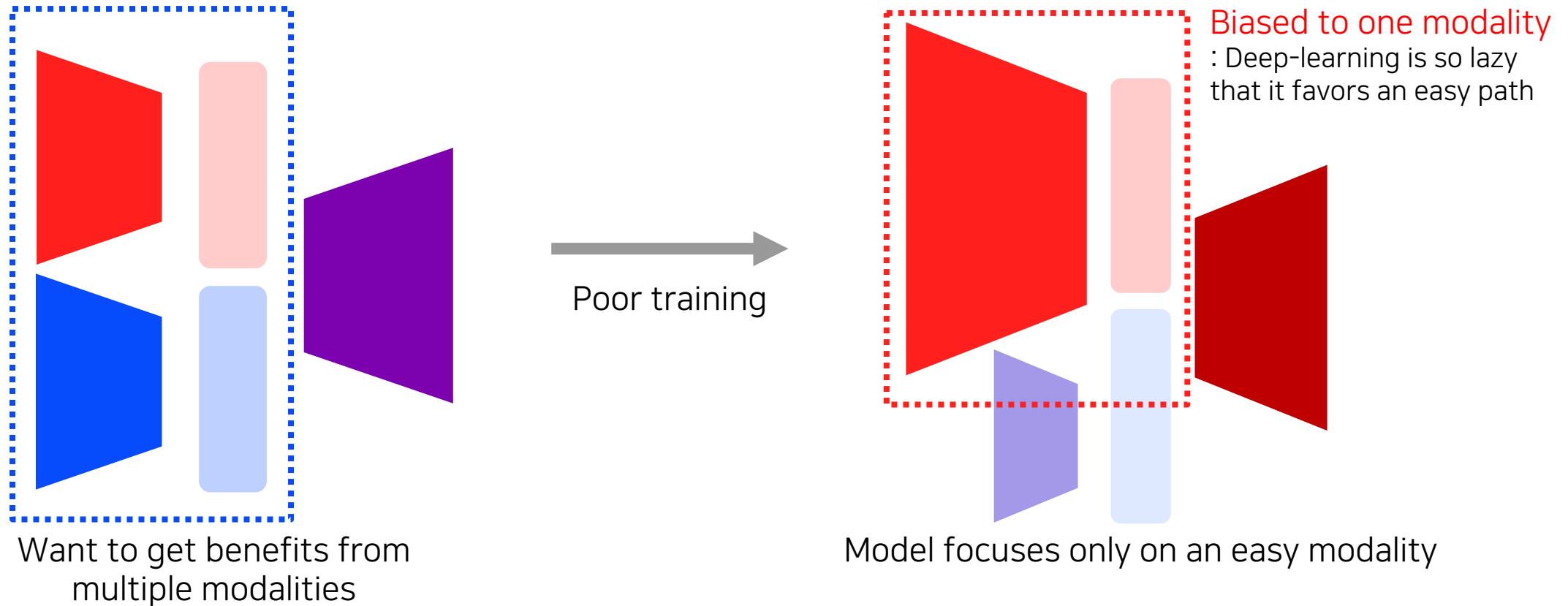


# 1. Overview of multi-modal learning

Multi-modal learning overview

Challenge (3) - May a model be biased on a specific modality

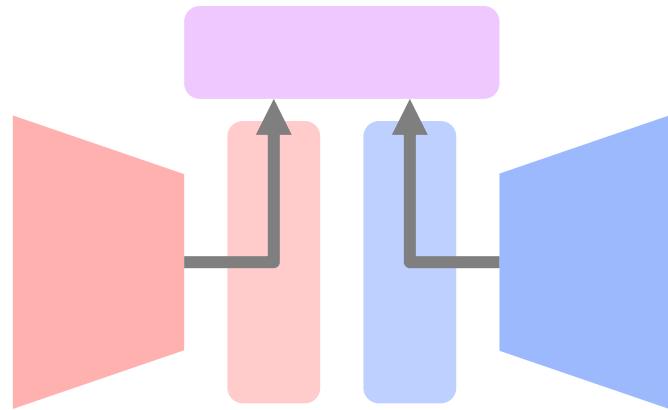
[Wang et al., CVPR 2020]



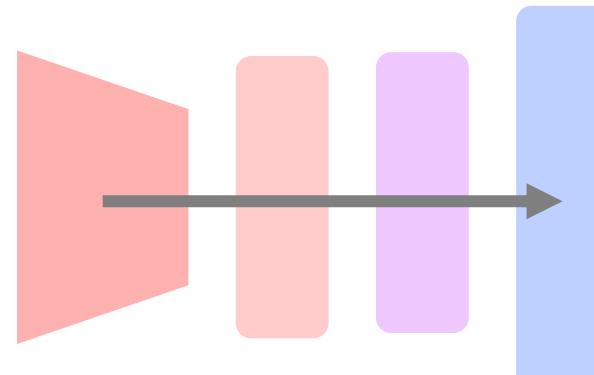
# 1. Overview of multi-modal learning

Multi-modal learning overview

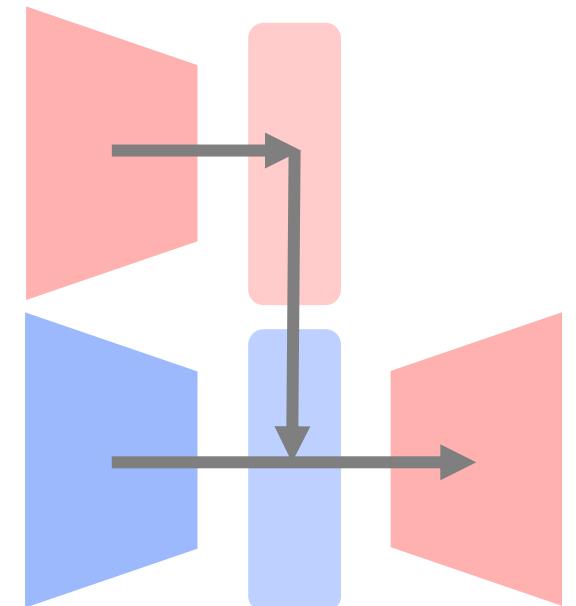
Despite the challenges, multi-modal learning is fruitful and important



Matching



Translating



Referencing

2.

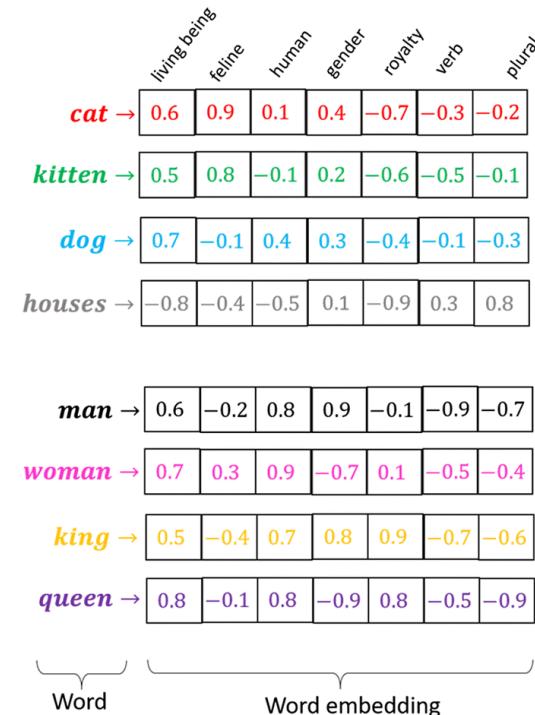
## Multi-modal tasks (1) - Visual data & Text

## 2.1 Text embedding

Multi-modal - Text

### Text embedding - Example

- Characters are hard to use in machine learning
- Map to dense vectors

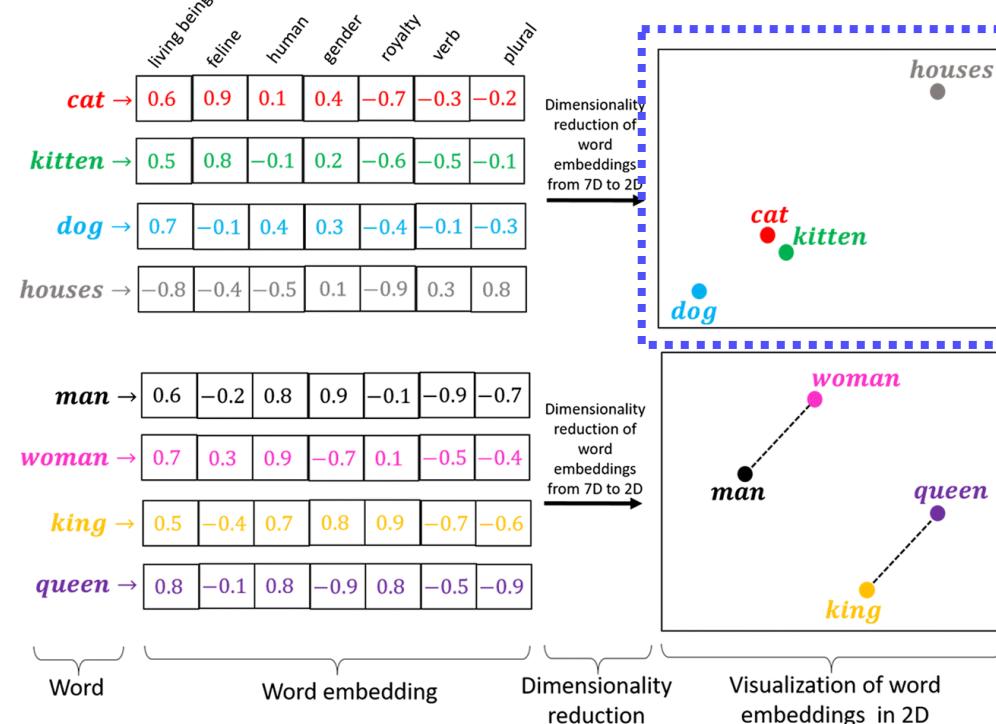


## 2.1 Text embedding

Multi-modal - Text

### Text embedding - Example

- Surprisingly, generalization power is obtained by learning dense representation
  - E.g., cat vs. kitten

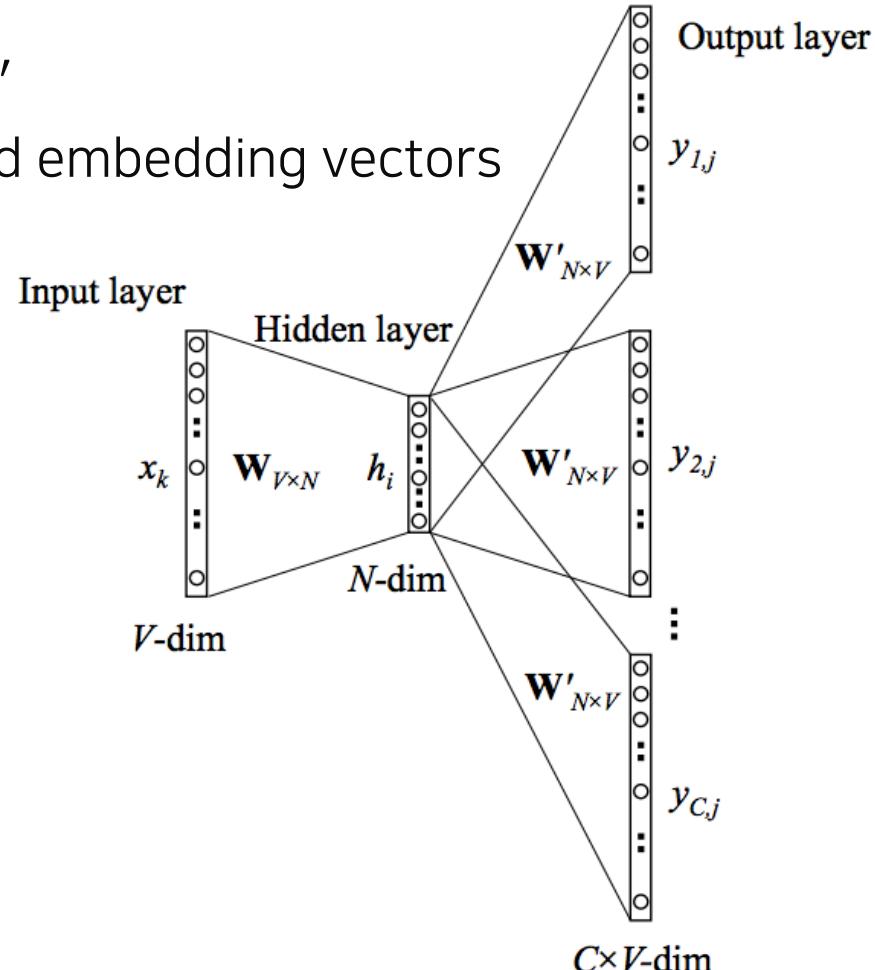


## 2.1 Text embedding

Multi-modal - Text

### word2vec - Skip-gram model

- Trained to learn  $W$  and  $W'$
- Rows in  $W$  represent word embedding vectors



## 2.1 Text embedding

Multi-modal - Text

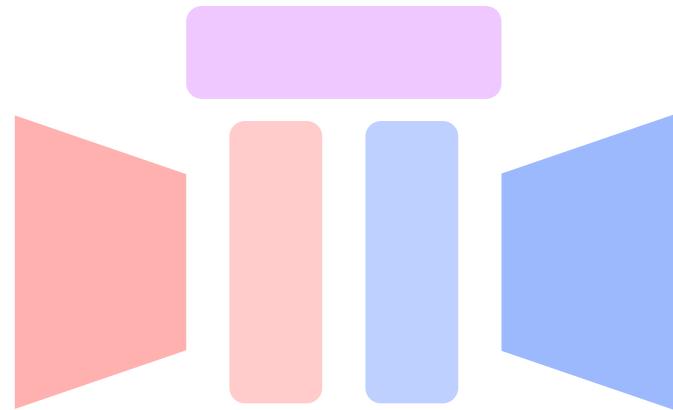
### word2vec - Skip-gram model

- Learning to predict neighboring  $N$  words for understanding relationships between words
  - Given a model with a window of size 5, the center words depend on 4 words

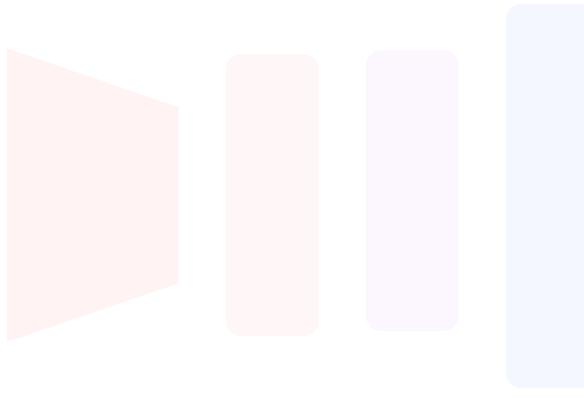
Source Text	Training Samples
The <span style="border: 1px solid black; padding: 2px;">quick</span> brown fox jumps over the lazy dog. ➔	(the, quick) (the, brown)
The quick <span style="border: 1px solid black; padding: 2px;">brown</span> fox jumps over the lazy dog. ➔	(quick, the) (quick, brown) (quick, fox)
The quick brown <span style="border: 1px solid black; padding: 2px;">fox</span> jumps over the lazy dog. ➔	(brown, the) (brown, quick) (brown, fox) (brown, jumps)
The quick brown fox <span style="border: 1px solid black; padding: 2px;">jumps</span> over the lazy dog. ➔	(fox, quick) (fox, brown) (fox, jumps) (fox, over)

## 2.2 Joint embedding

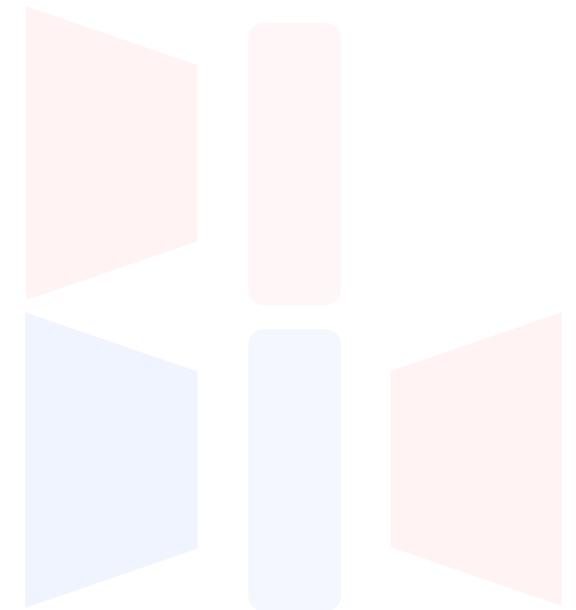
Multi-modal - Text



Matching



Translating



Referencing

## 2.2 Joint embedding

Multi-modal - Text

### Application – Image tagging

[Srivastava and Slakhutdinov, JMLR 2014]

- Can generate tags of a given image, and retrieve images by a tag keyword as well



Tag generation of an image

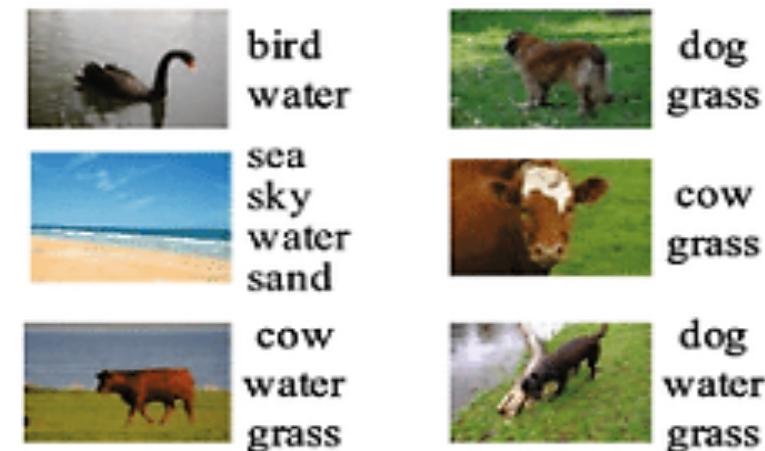


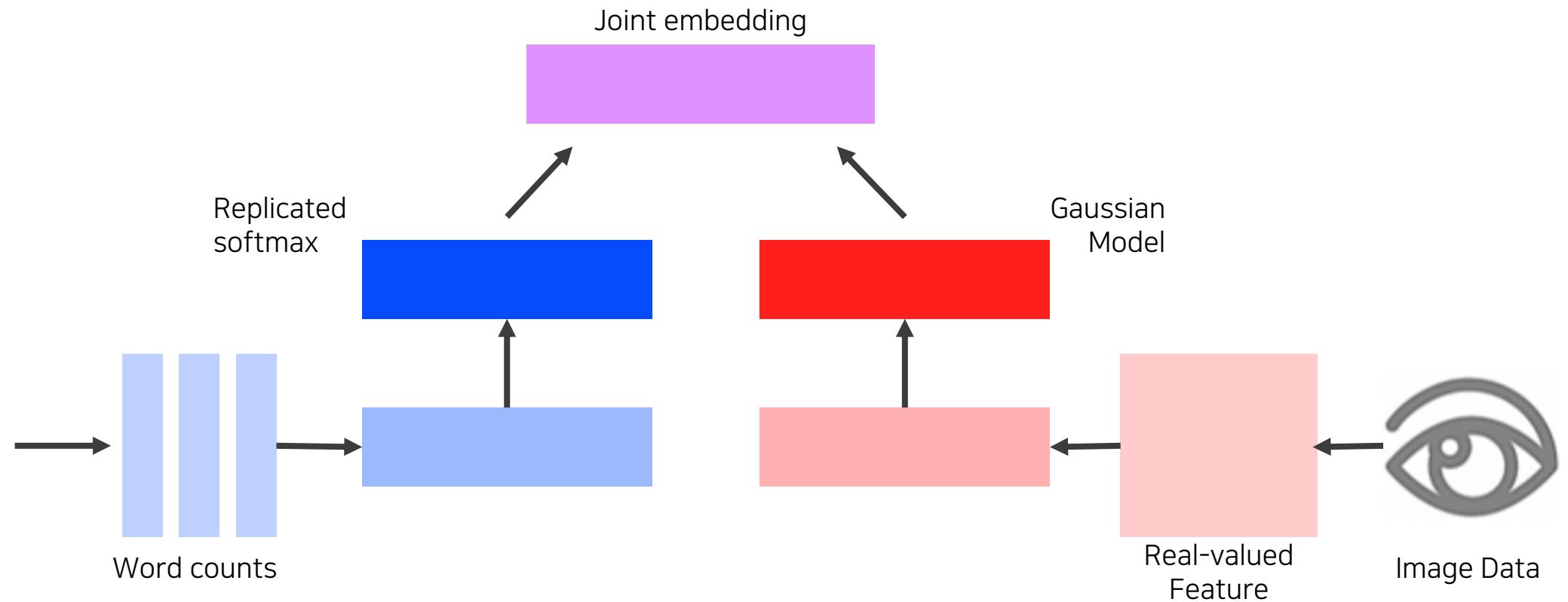
Image retrieval by a tag keyword

## 2.2 Joint embedding

Multi-modal - Text

Image tagging - Combining pre-trained unimodal models

[Srivastava and Slakhutdinov, JMLR 2014]

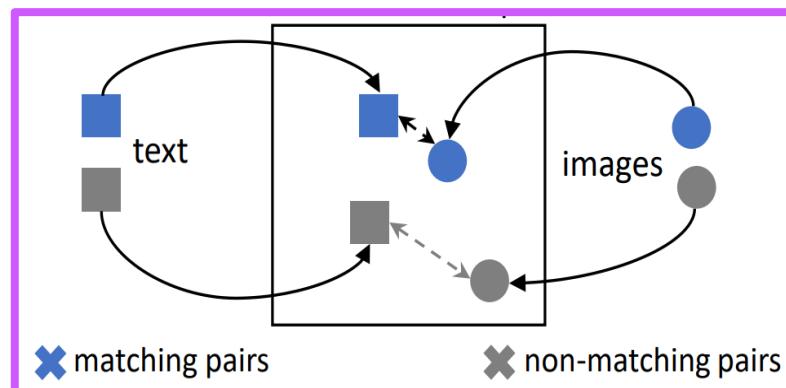


## 2.2 Joint embedding

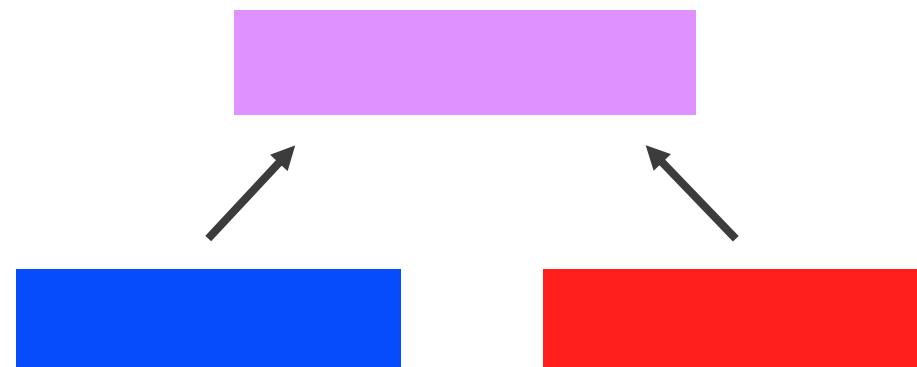
Multi-modal - Text

Image tagging - Metric learning in visual-semantic space

[Srivastava and Slakhutdinov, JMLR 2014]



Joint visual-semantic embedding space



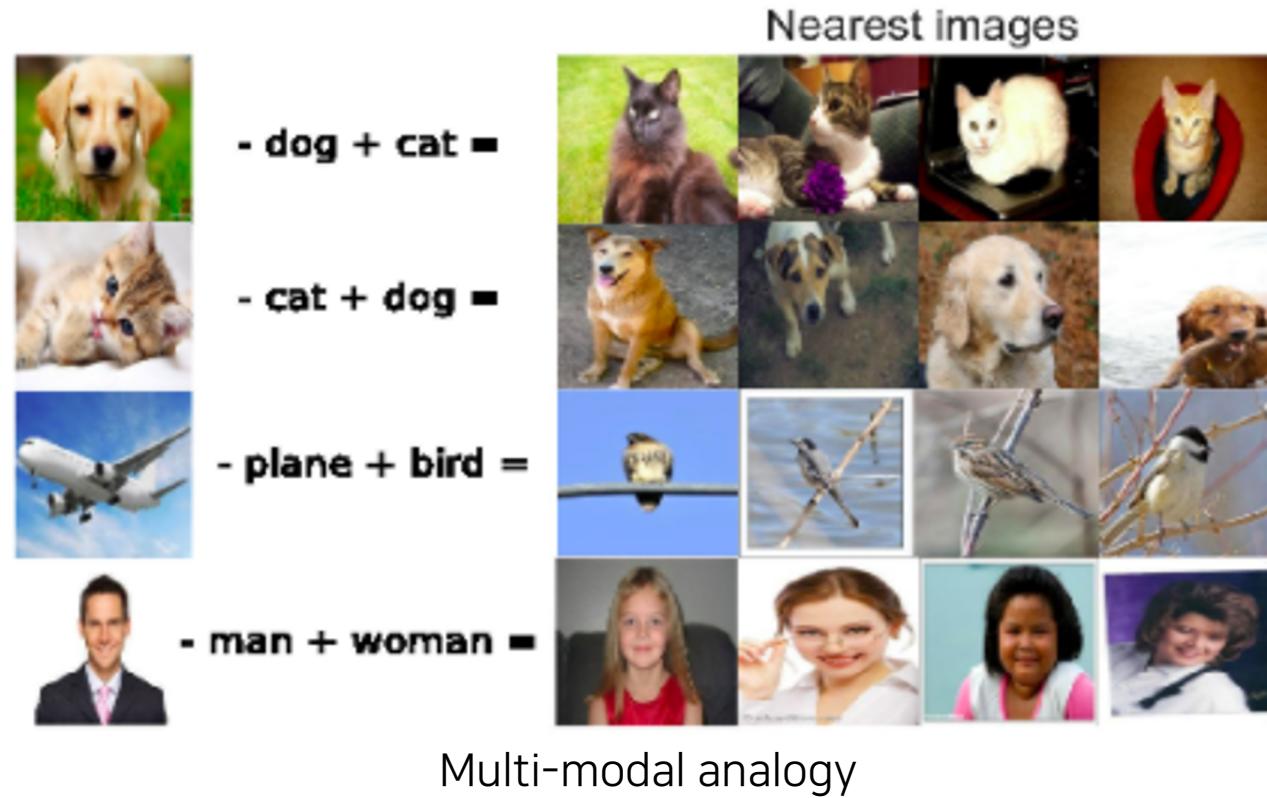
## 2.2 Joint embedding

Multi-modal - Text

Image tagging - Interesting property

[Srivastava and Slakhutdinov, JMLR 2014]

- Surprisingly, learned embeddings hold analogy relationships between visual and text data



## 2.2 Joint embedding

Multi-modal - Text

Application – Image & food recipe retrieval

[Marin et al., TPAMI 2019]

Query Image



True ingr.

whole milk  
half - and - half cr  
white sugar  
lemon extract  
ground cinnamon  
frozen blueberries  
vanilla wafers  
ice cubes

Retrieved ingr.

berries  
strawberry yogurt  
banana  
milk  
white sugar



butter  
garlic cloves  
all - purpose flour  
kosher salt  
milk  
chicken broth  
mozzarella cheese  
parmesan cheese  
onion

1 box any pasta you  
ground beef  
1 envelope taco seas  
water  
1/2 packages cream c  
cheese



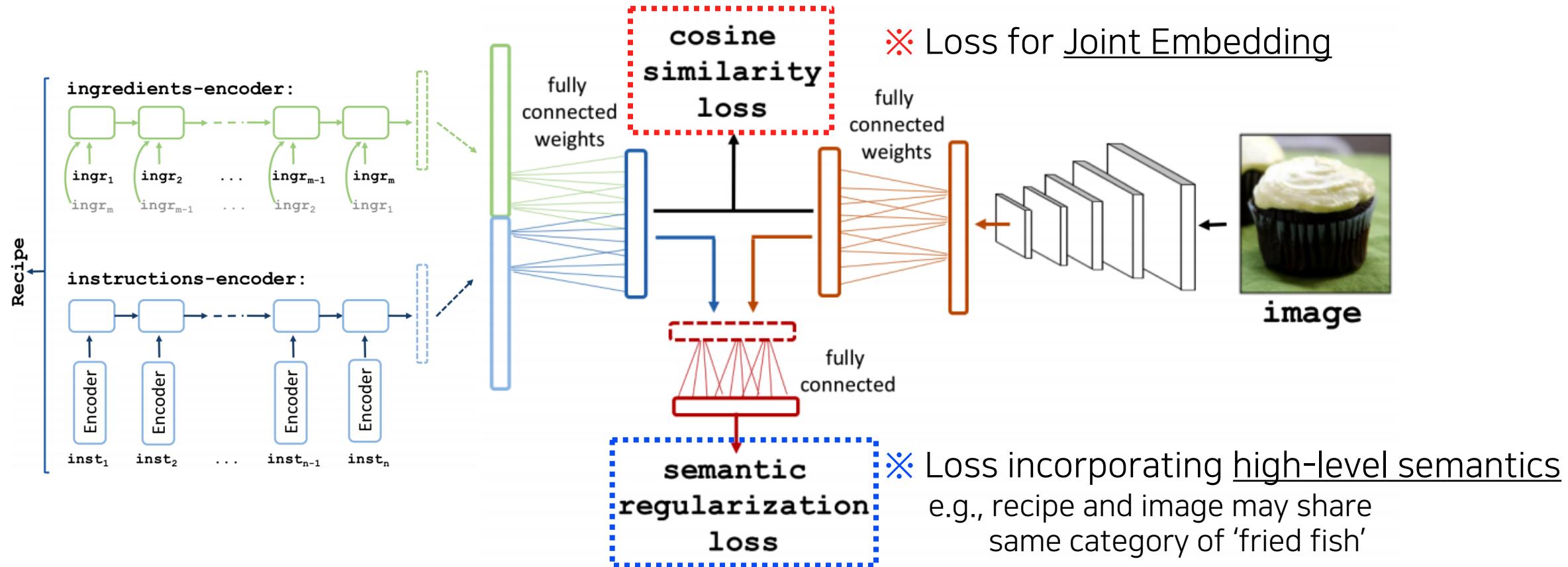
Retrieved ingredients and images given the query images

## 2.2 Joint embedding

Multi-modal - Text

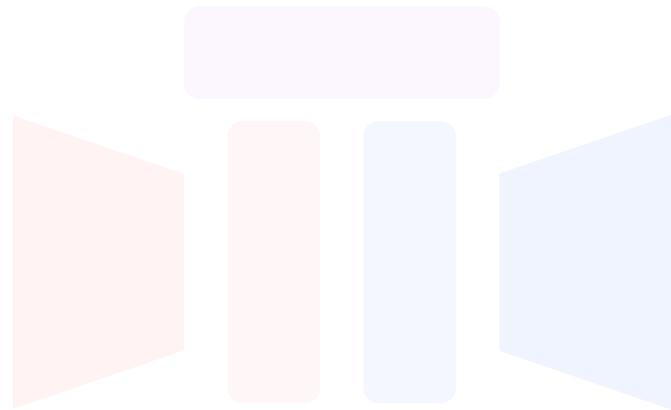
Recipe text (sentence) vs. food image

[Marin et al., TPAMI 2019]

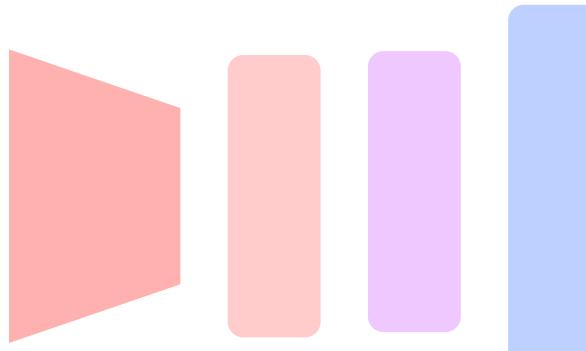


## 2.3 Cross modal translation

Multi-modal - Text



Matching



Translating



Referencing

## 2.3 Cross modal translation

Multi-modal - Text

### Application – Image captioning



a cat sitting on a  
suitcase on the floor



a cat is sitting on  
a tree branch



a dog is running in the  
grass with a frisbee



a white teddy bear  
sitting in the grass



two people walking on  
the beach with  
surfboards



a tennis player in  
action on the court



two giraffes standing  
in a grassy field



a man riding a dirt  
bike on a dirt track

## 2.3 Cross modal translation

Multi-modal - Text

Captioning as image-to-sentence - CNN for image & RNN for sentence

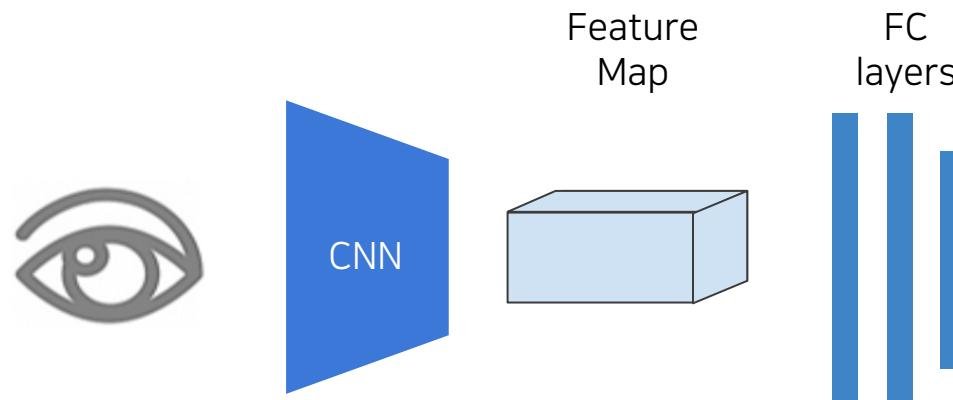
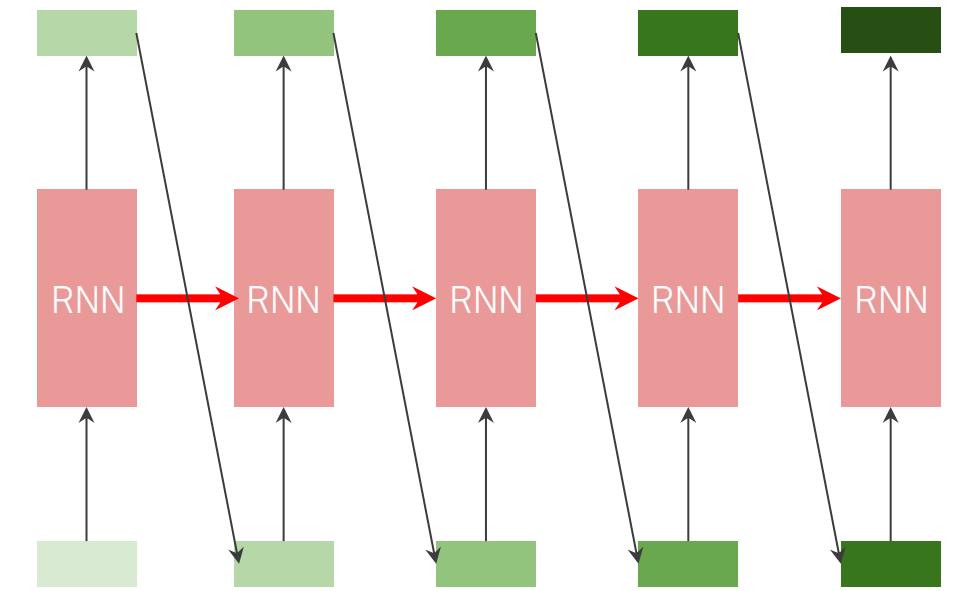


Image representation by CNN



Sentence prediction by RNN

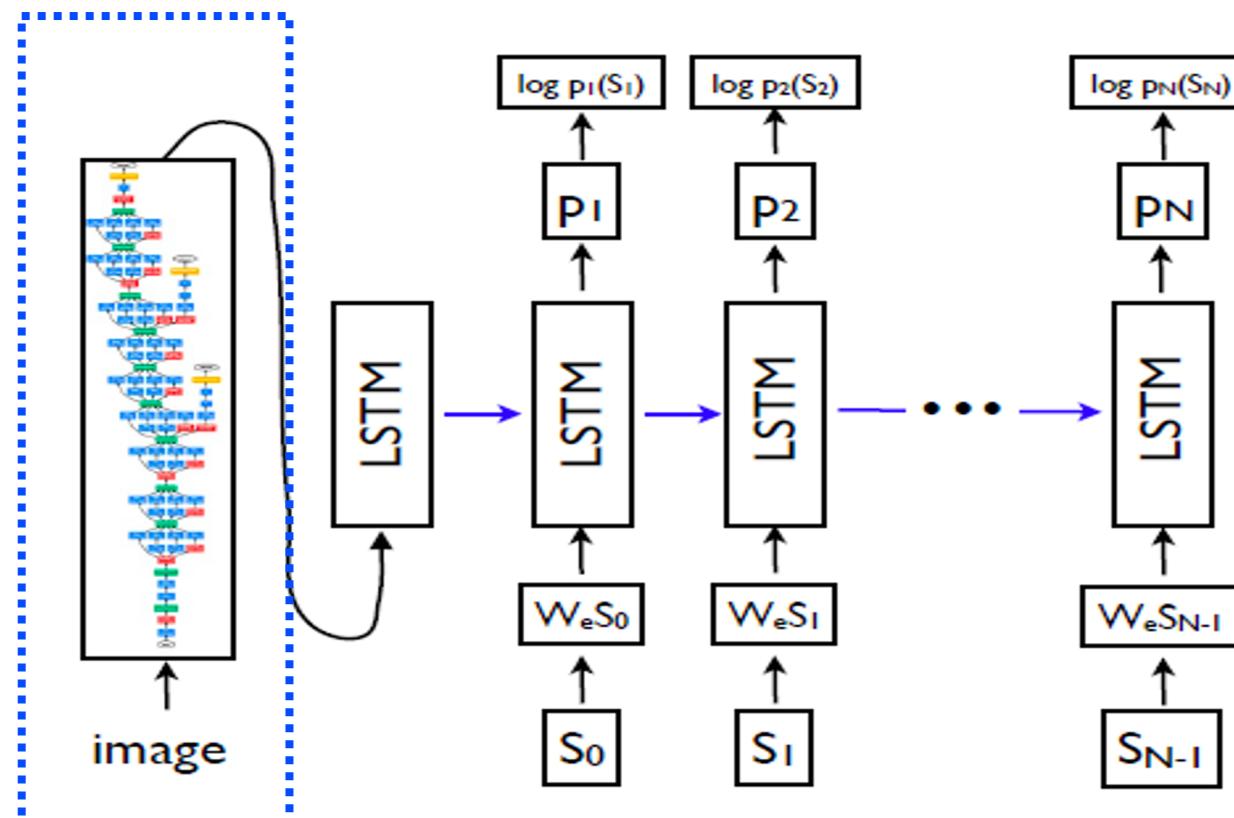
## 2.3 Cross modal translation

Multi-modal - Text

Show and tell

[Vinyals et al., CVPR 2015]

- Encoder: CNN model pre-trained on ImageNet
- Decoder: LSTM module



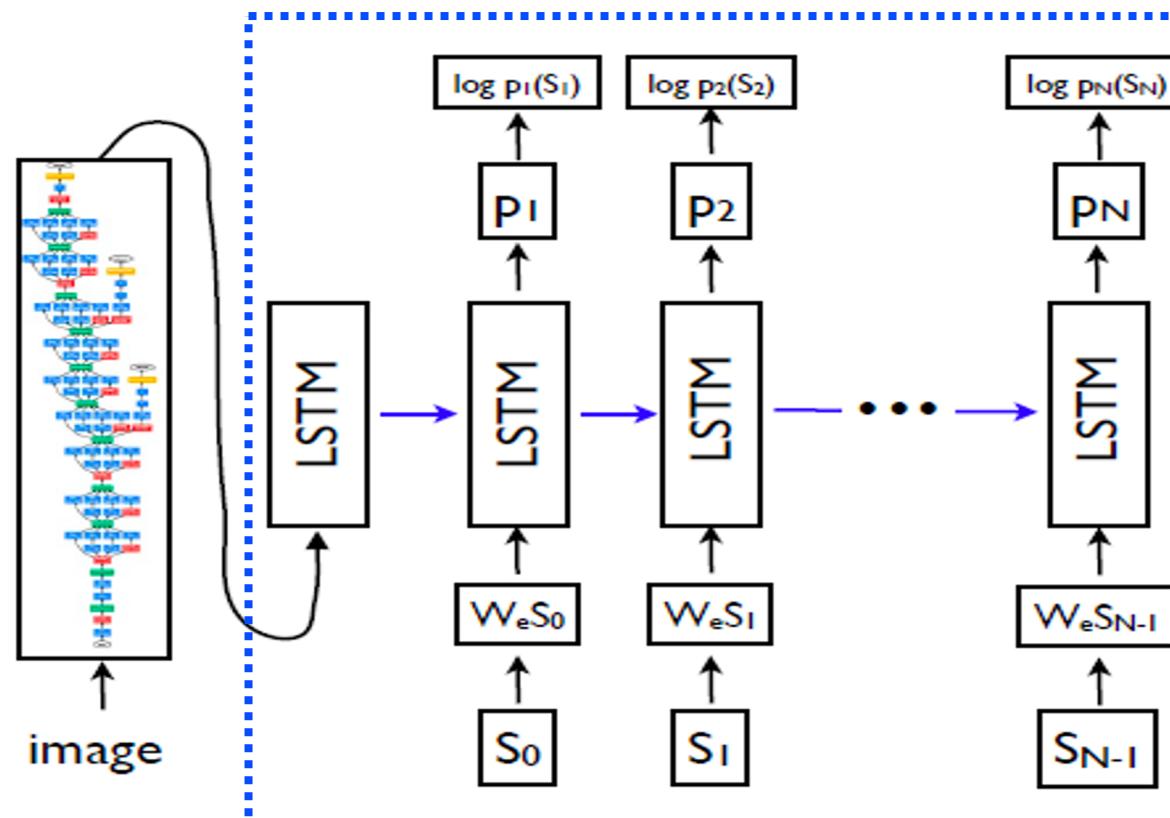
## 2.3 Cross modal translation

Multi-modal - Text

Show and tell

[Vinyals et al., CVPR 2015]

- Encoder: CNN model pre-trained on ImageNet
- Decoder: LSTM module



## 2.3 Cross modal translation

Multi-modal - Text

Show, attend, and tell - Example results

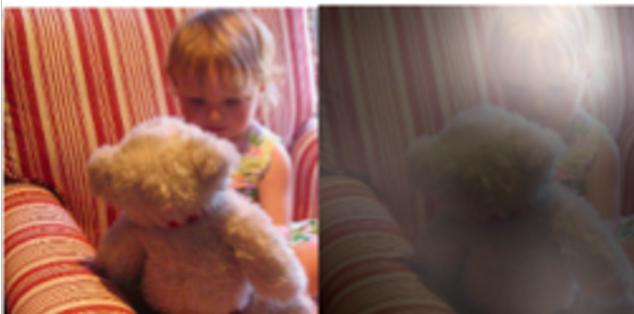
[Xu et al., ICML 2015]



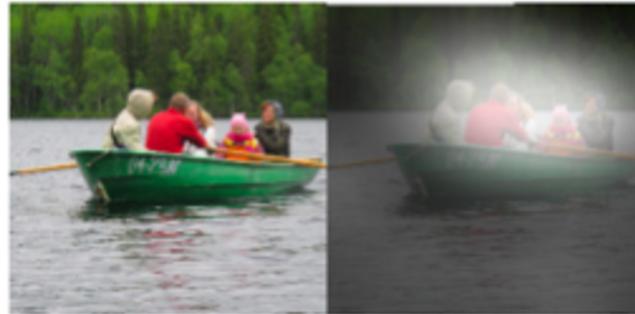
A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.



A little girl sitting on a bed with a teddy bear.



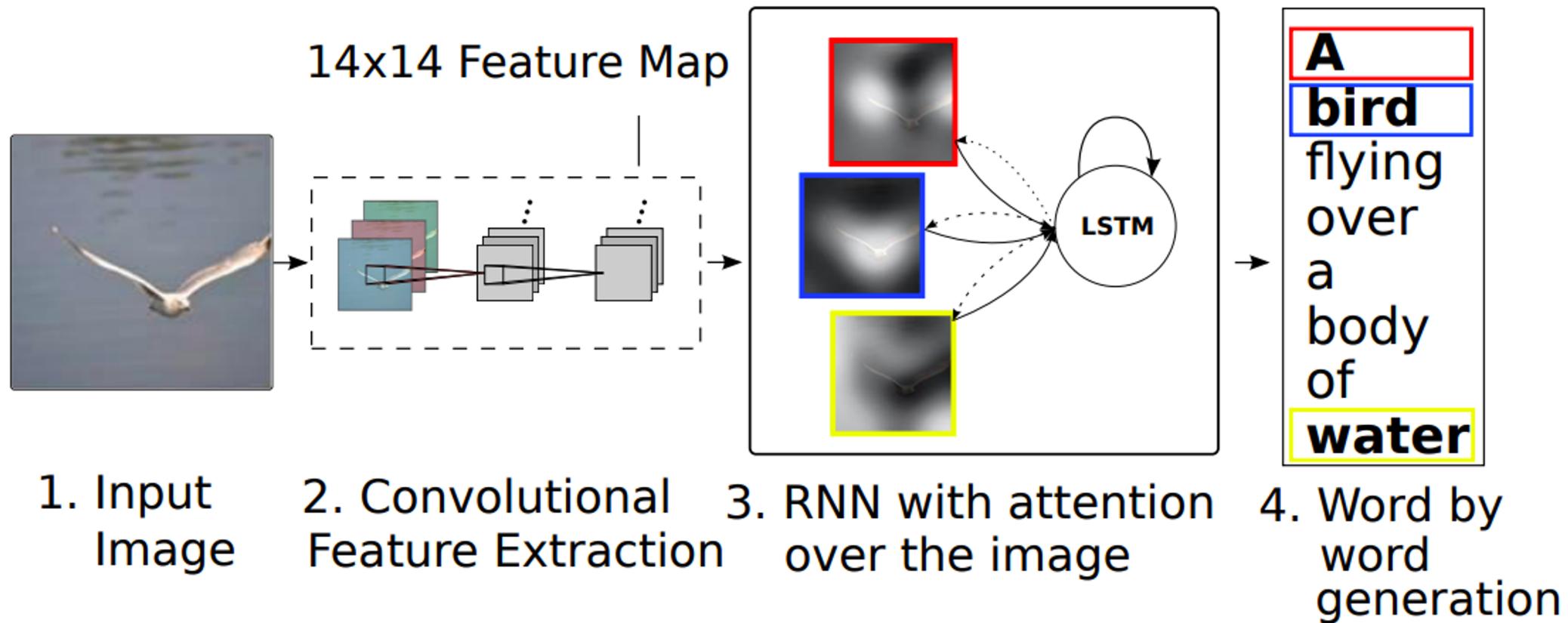
A group of people sitting on a boat in the water.

## 2.3 Cross modal translation

Multi-modal - Text

Show, attend, and tell

[Xu et al., ICML 2015]



## 2.3 Cross modal translation

Multi-modal - Text

Show, attend, and tell - Attention

[Xu et al., ICML 2015]

CC BY-SA 2.0 licensed image



Human attention shifts on face image

## 2.3 Cross modal translation

Multi-modal - Text

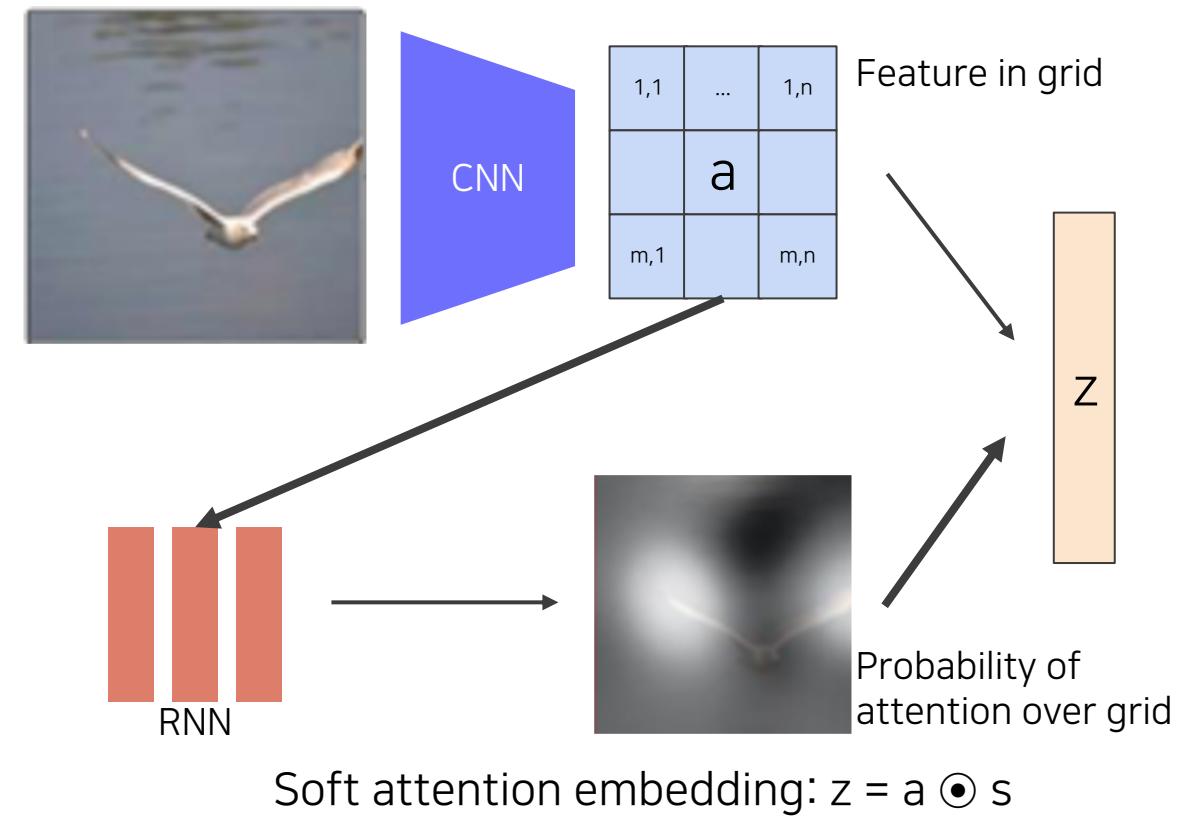
Show, attend, and tell - Soft Attention

[Xu et al., ICML 2015]

CC BY-SA 2.0 licensed image



Human attention shifts on face image

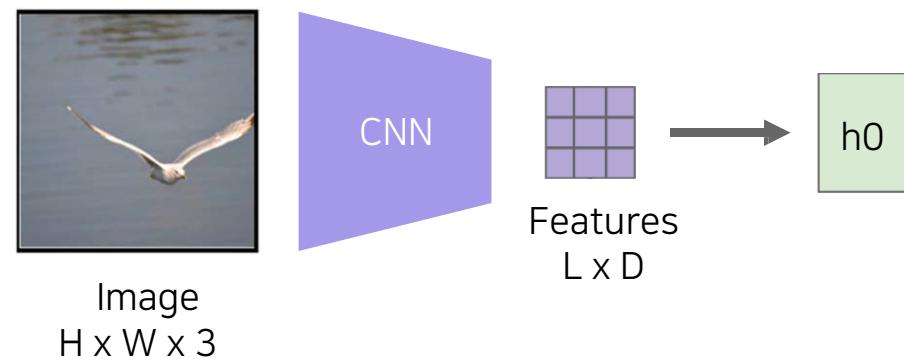


## 2.3 Cross modal translation

Multi-modal - Text

Show, attend, and tell - Inference

[Xu et al., ICML 2015]

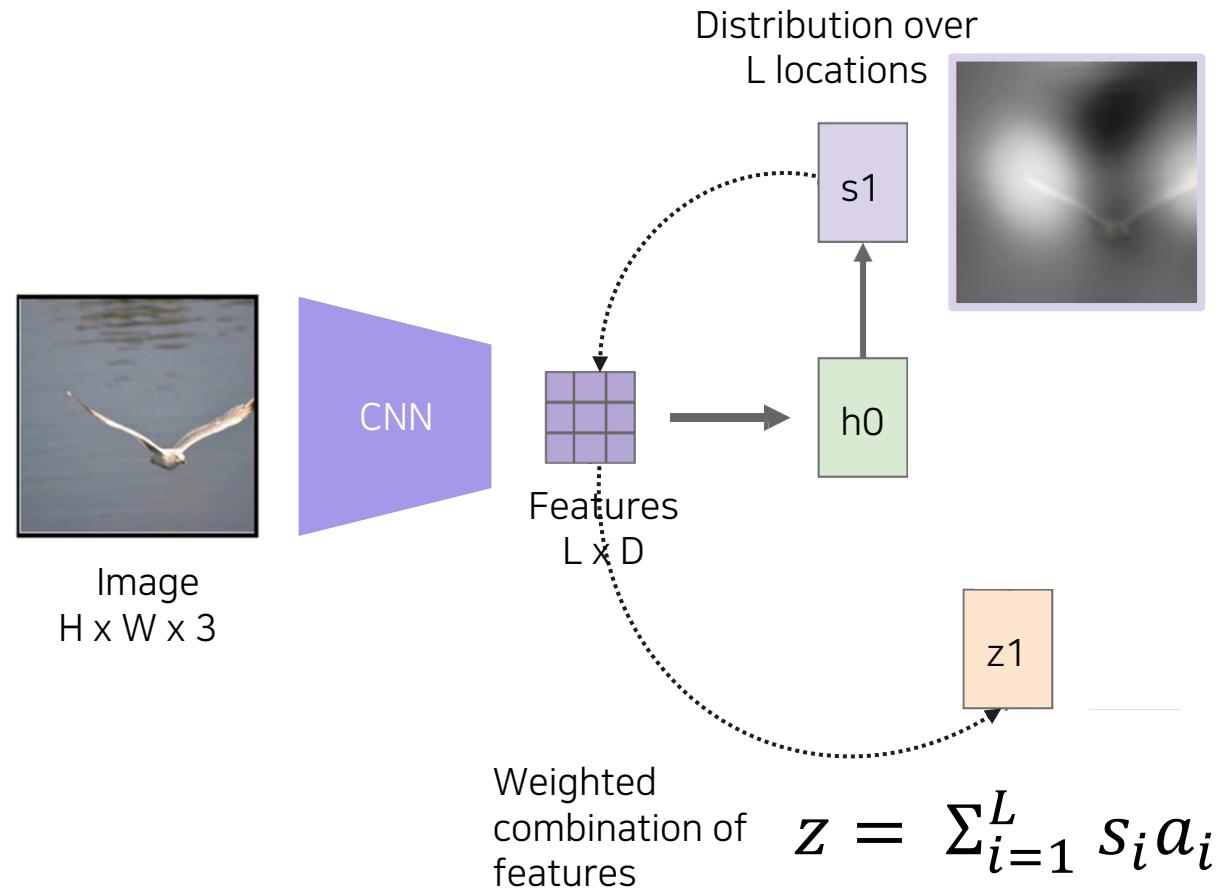


## 2.3 Cross modal translation

Multi-modal - Text

Show, attend, and tell - Inference

[Xu et al., ICML 2015]

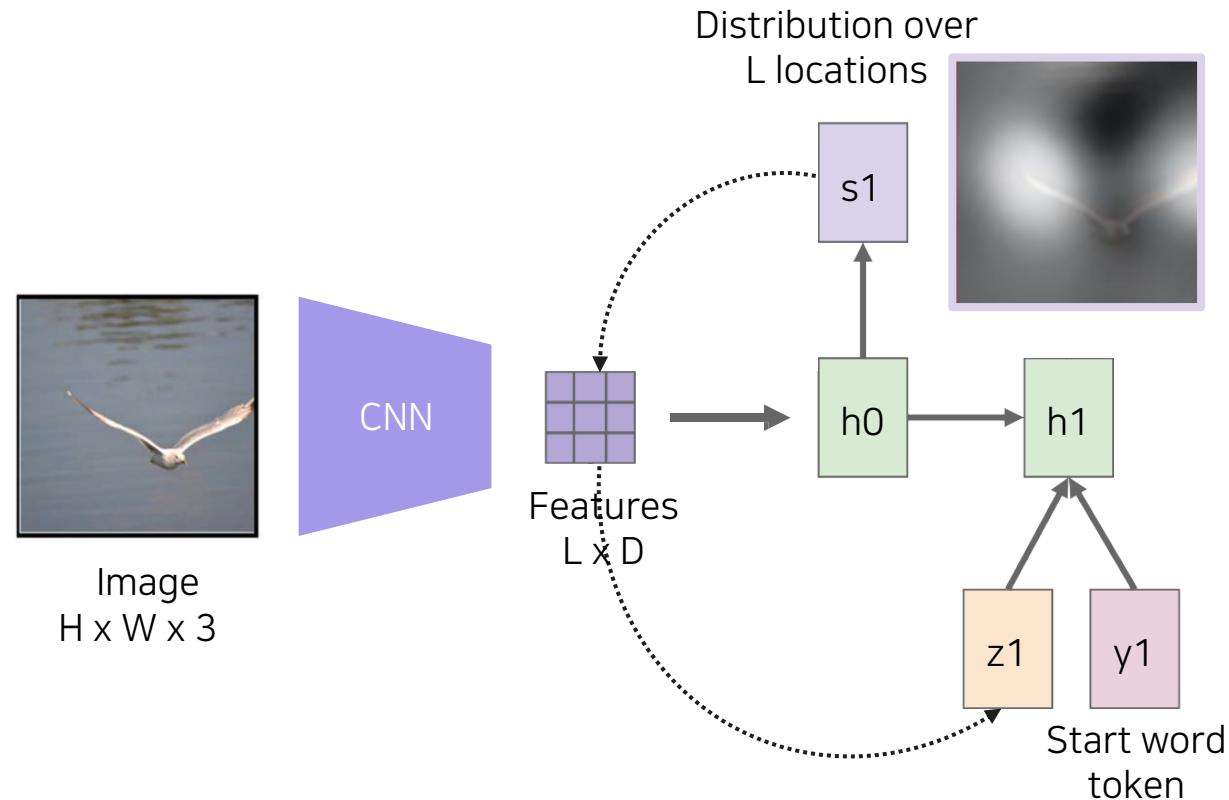


## 2.3 Cross modal translation

Multi-modal - Text

Show, attend, and tell - Inference

[Xu et al., ICML 2015]

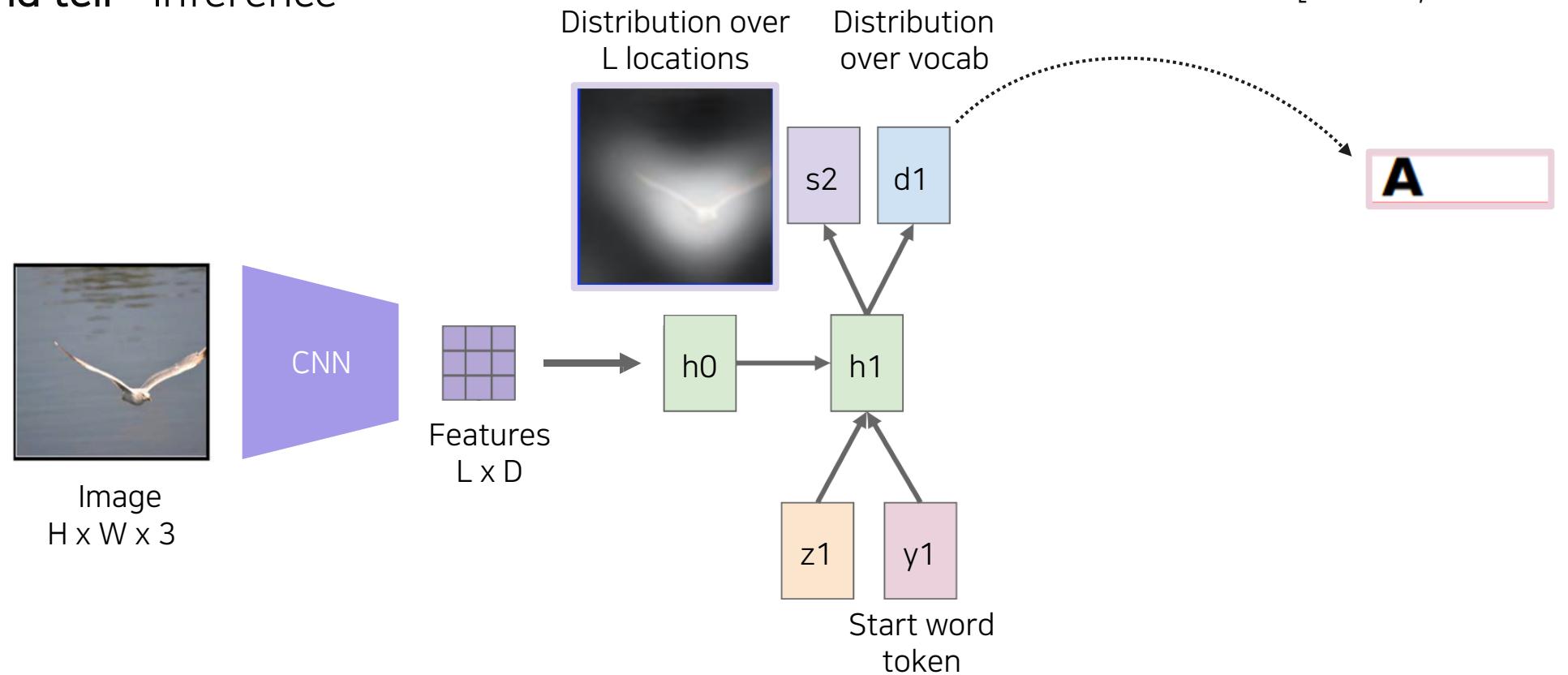


## 2.3 Cross modal translation

Multi-modal - Text

Show, attend, and tell - Inference

[Xu et al., ICML 2015]

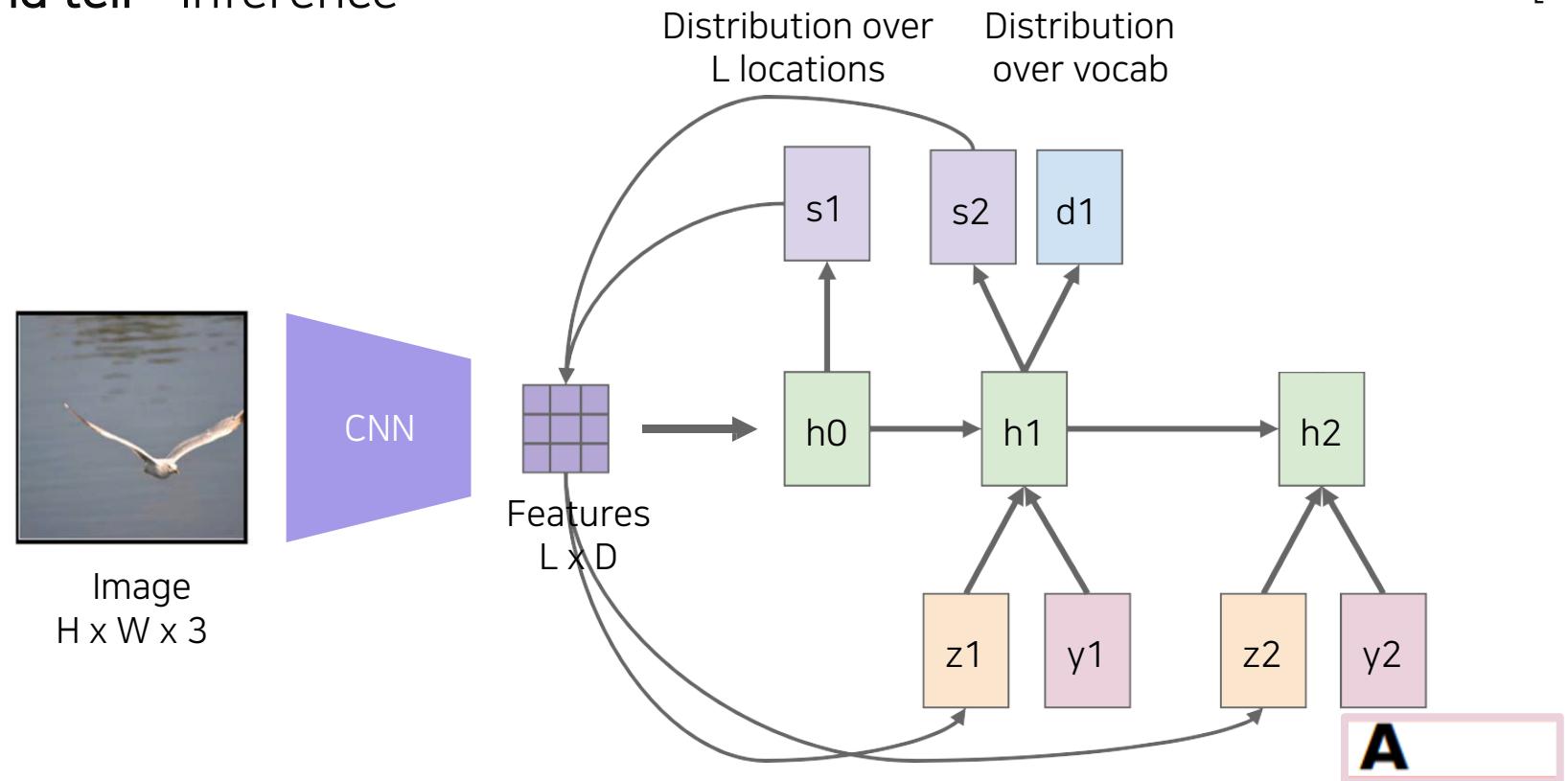


## 2.3 Cross modal translation

Multi-modal - Text

Show, attend, and tell - Inference

[Xu et al., ICML 2015]

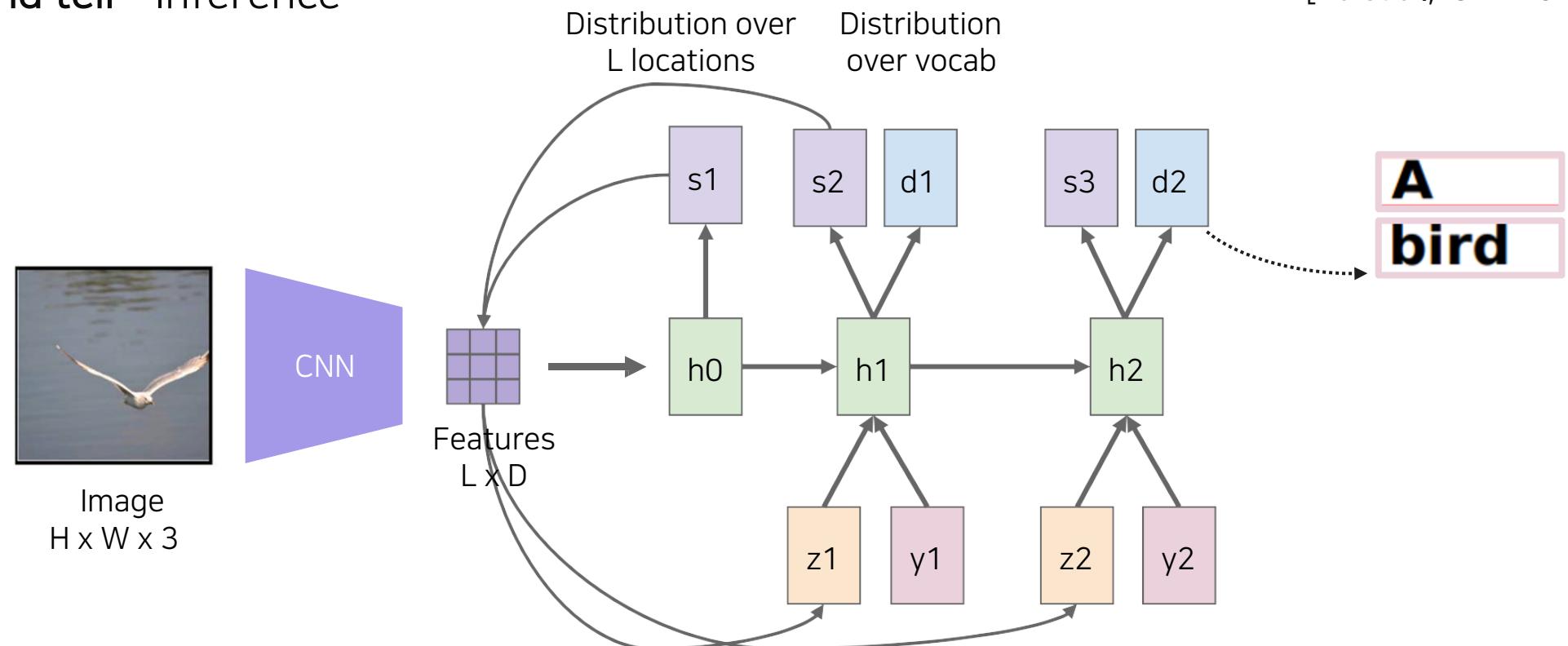


## 2.3 Cross modal translation

Multi-modal - Text

Show, attend, and tell - Inference

[Xu et al., ICML 2015]

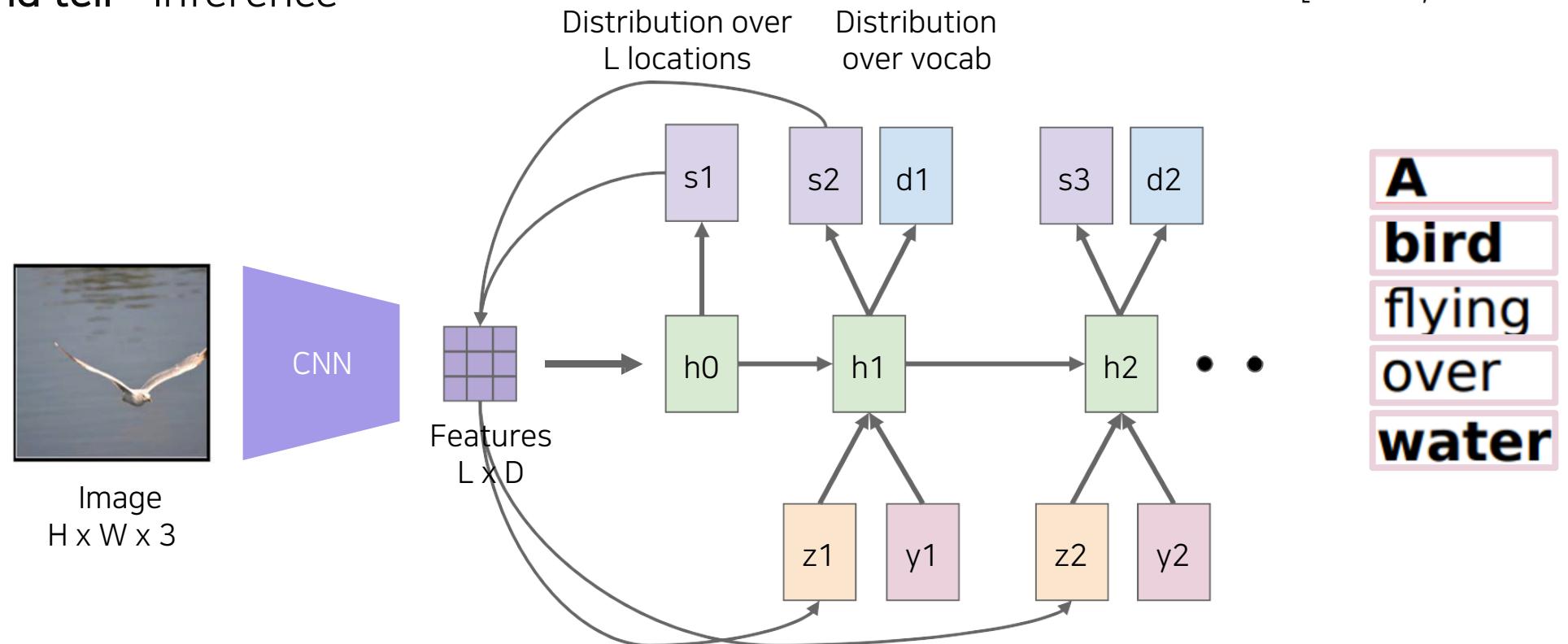


## 2.3 Cross modal translation

Multi-modal - Text

Show, attend, and tell - Inference

[Xu et al., ICML 2015]



## 2.3 Cross modal translation

Multi-modal - Text

Text-to-image by generative model - Example

[Reed et al., ICML 2016]

this small bird has a pink breast and crown, and black primaries and secondaries.



this magnificent fellow is almost all black with a red crest, and white cheek patch.



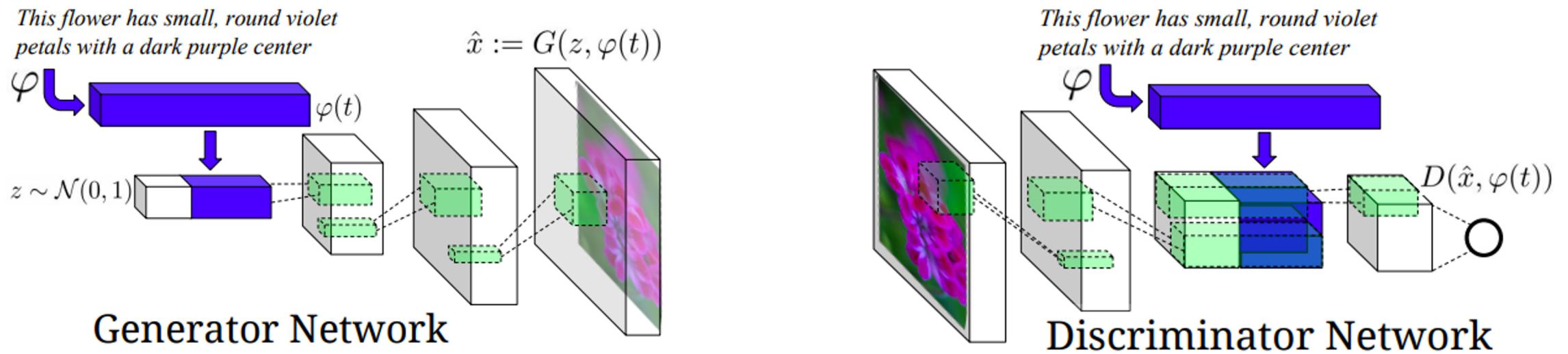
(Top) Query sentences and (Bottom) corresponding 6 generated images

## 2.3 Cross modal translation

Multi-modal - Text

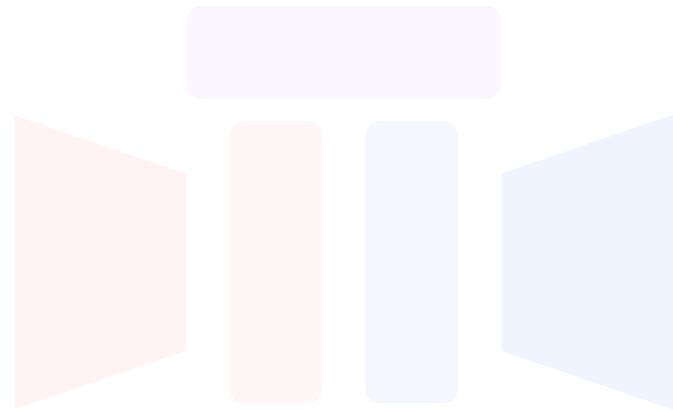
Text-to-image by generative model

[Reed et al., ICML 2016]

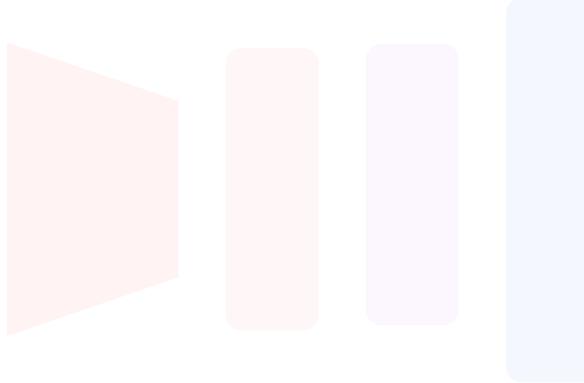


## 2.4 Cross modal reasoning

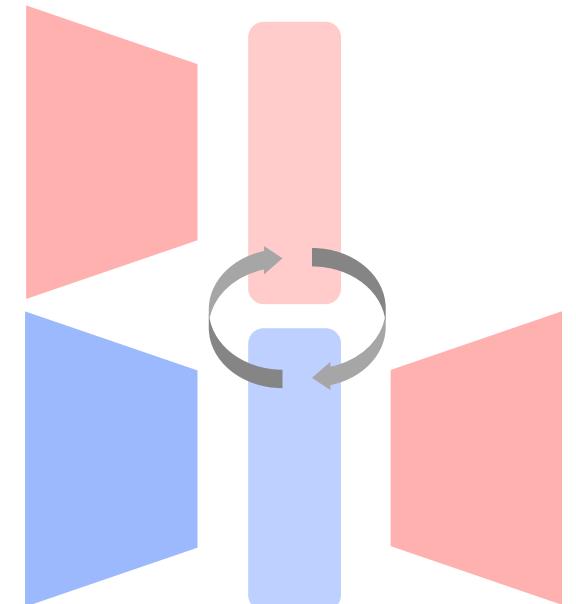
Multi-modal - Text



Matching



Translating



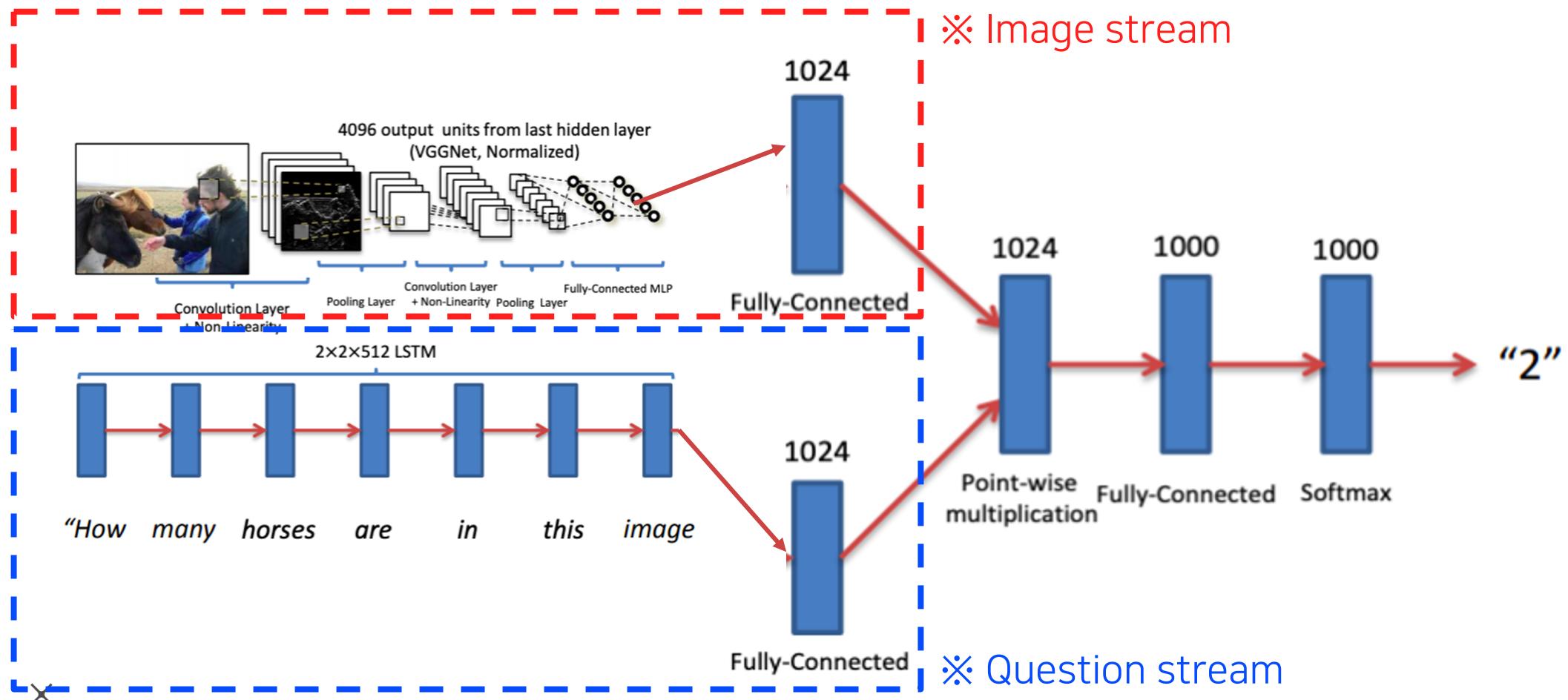
Referencing

## 2.4 Cross modal reasoning

Multi-modal - Text

Visual question answering - Multiple streams

[Antol et al., ICCV 2015]

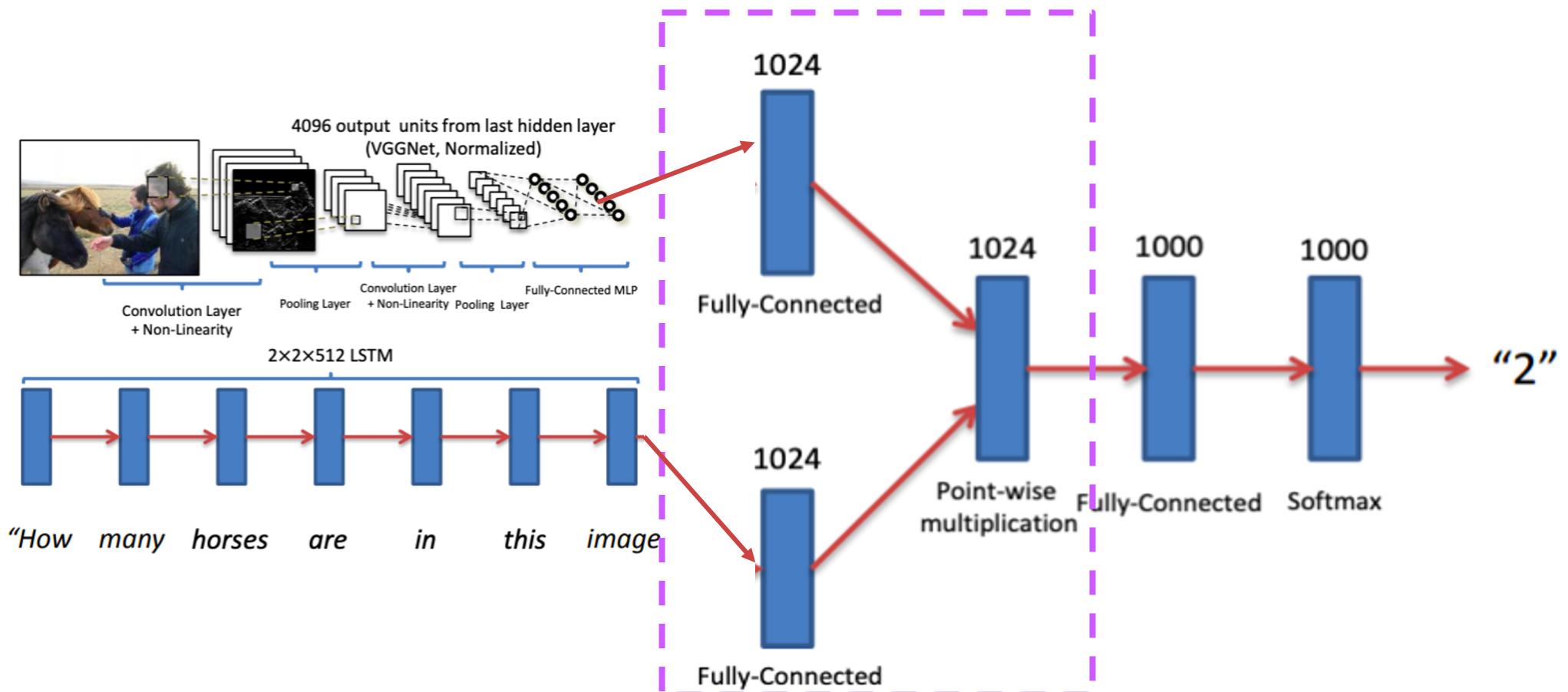


## 2.4 Cross modal reasoning

Multi-modal - Text

Visual question answering – Joint embedding

[Antol et al., ICCV 2015]

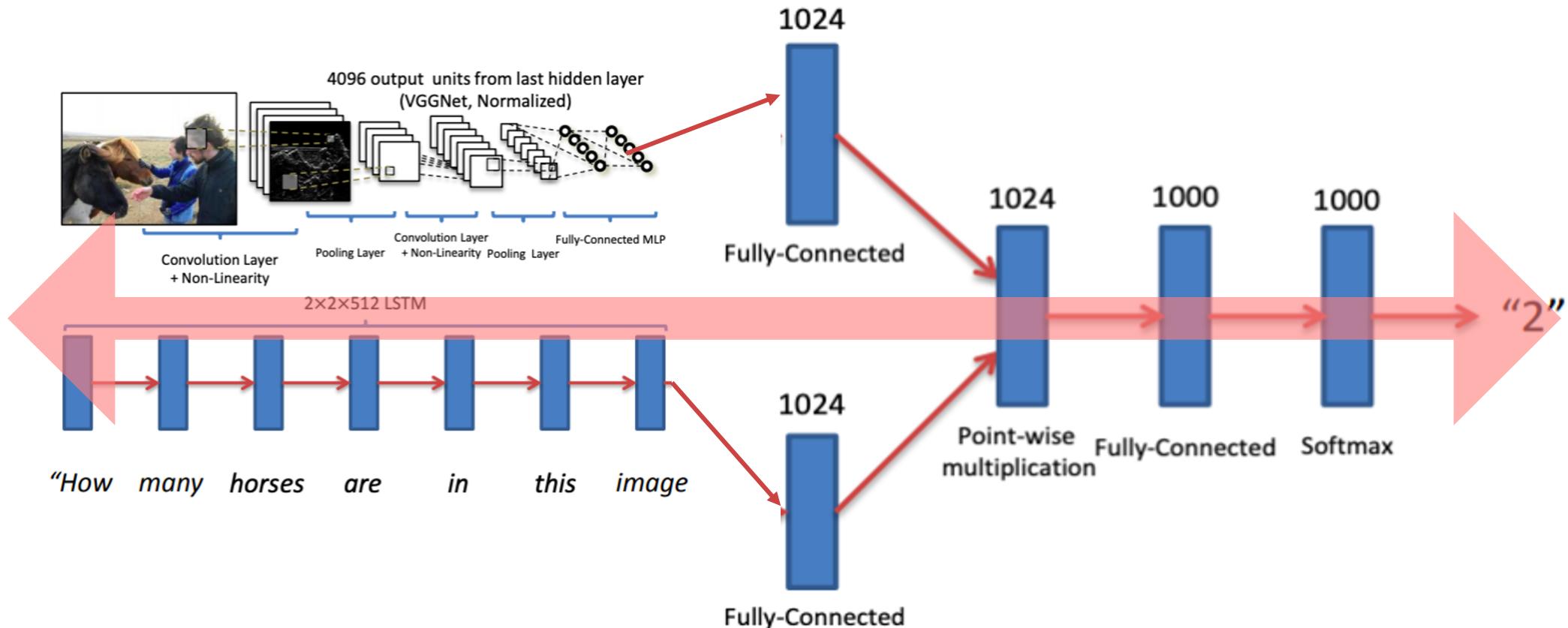


## 2.4 Cross modal reasoning

Multi-modal - Text

Visual question answering - End-to-end training

[Antol et al., ICCV 2015]



3.

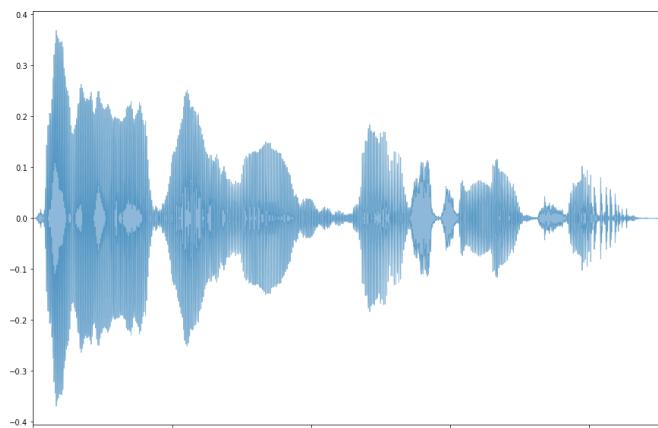
## Multi-modal tasks (2) - Visual data & Audio

# 3.1 Sound representation

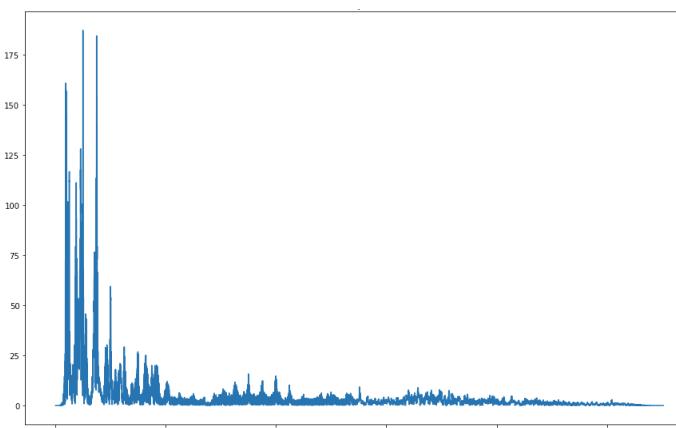
Multi-modal - Audio

## Sound representation

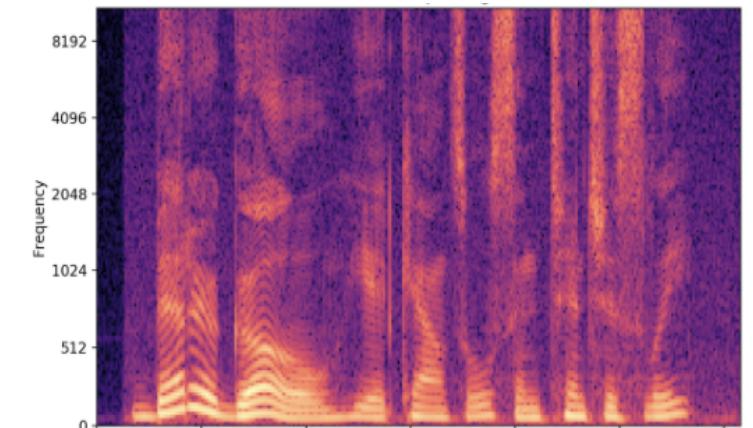
- Acoustic feature extraction from waveform to spectrogram



Waveform



Power spectrum



Spectrogram

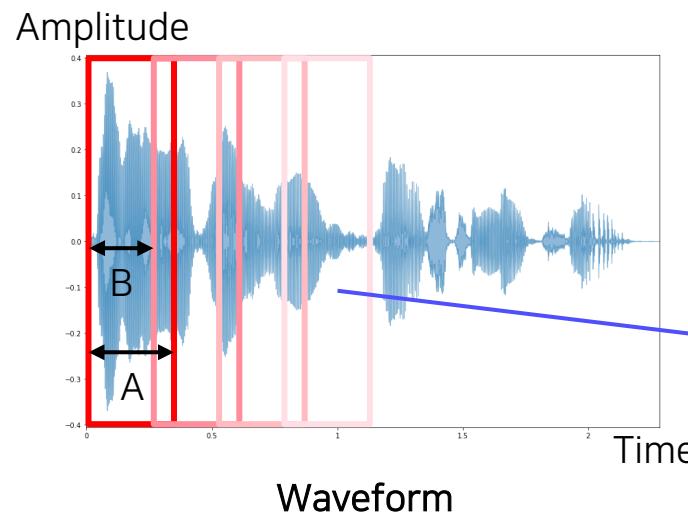


# 3.1 Sound representation

Multi-modal - Audio

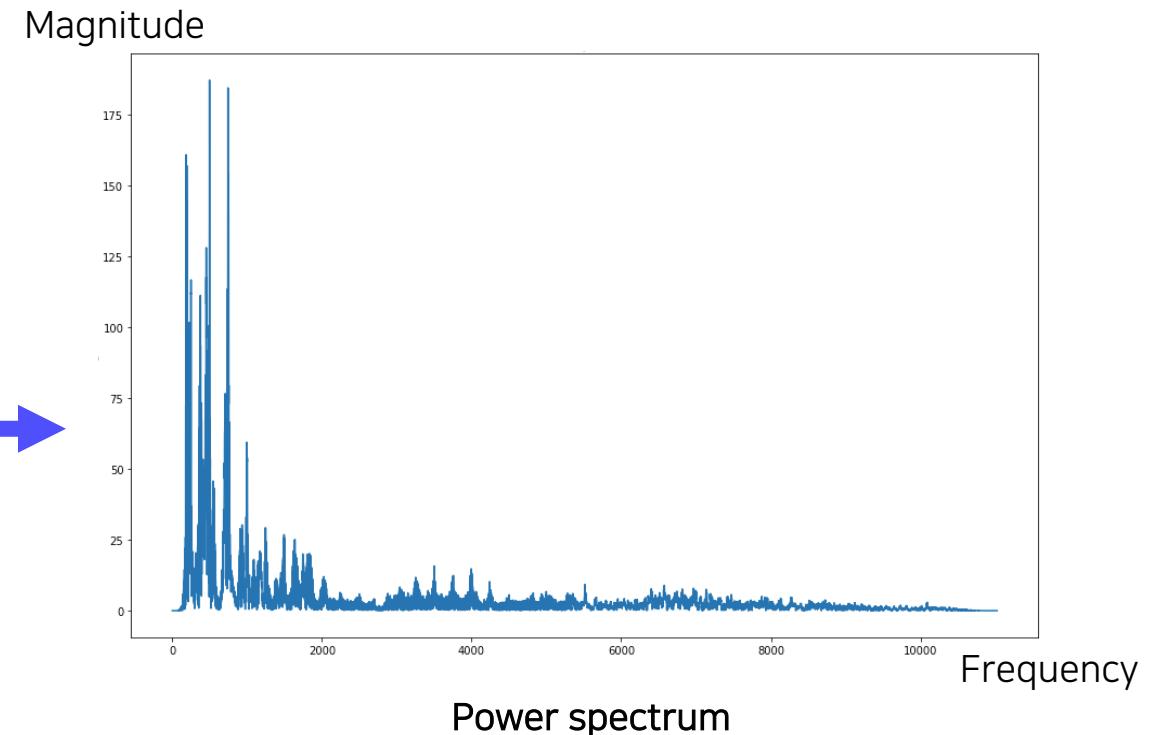
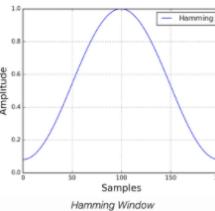
## Sound representation - Fourier transform

- Short-time Fourier transform (STFT)  
: Fourier transform (FT) on windowed waveform results in frequency-magnitude graph



A : 20~25ms  
B : 10ms

Hamming window

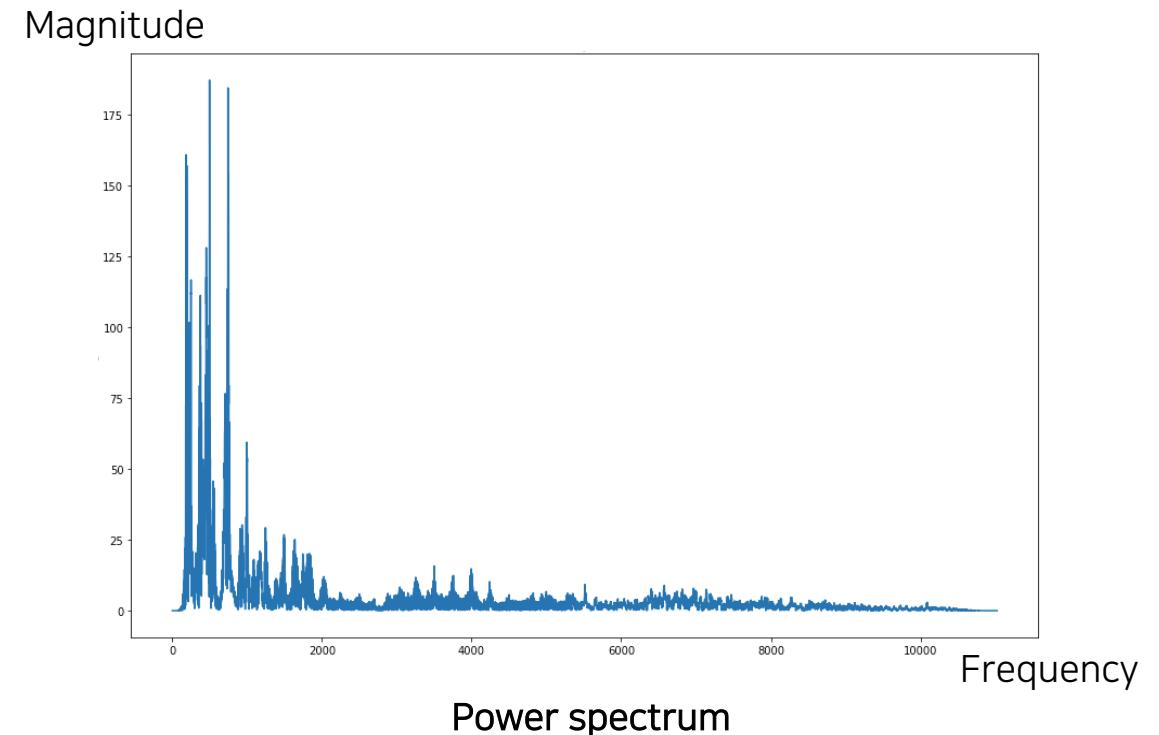
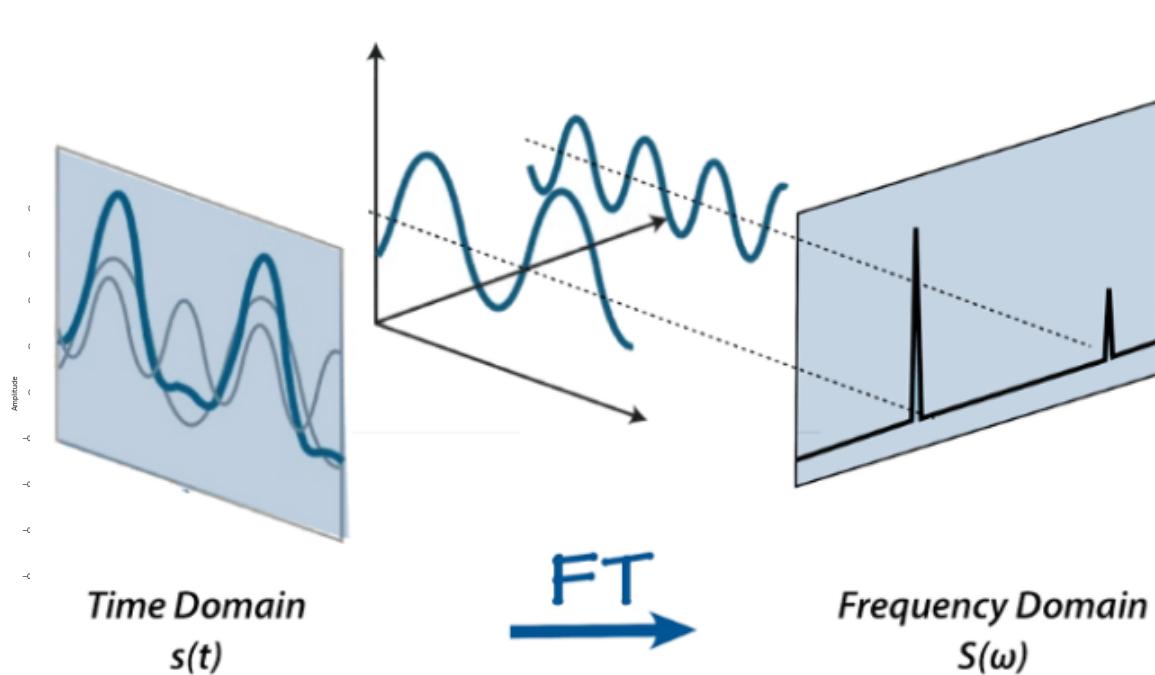


## 3.1 Sound representation

Multi-modal - Audio

### Sound representation - Fourier transform

- FT decomposes an input signal into constituent frequencies

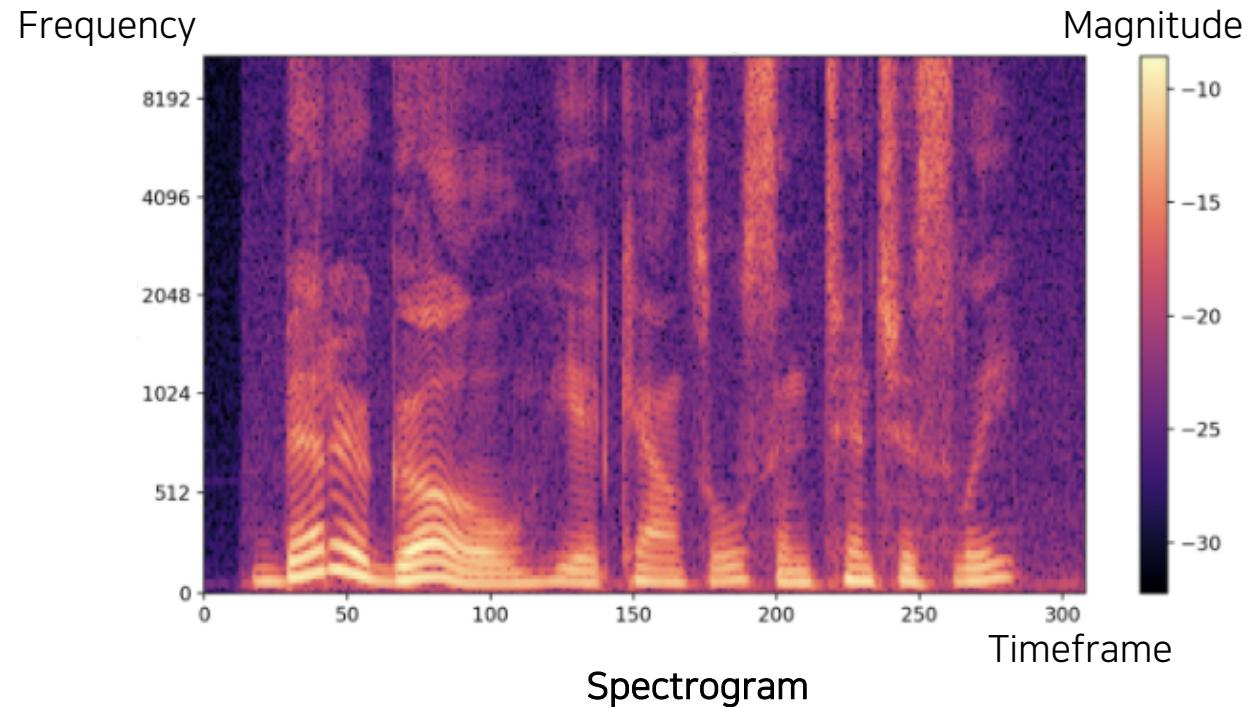


## 3.1 Sound representation

Multi-modal - Audio

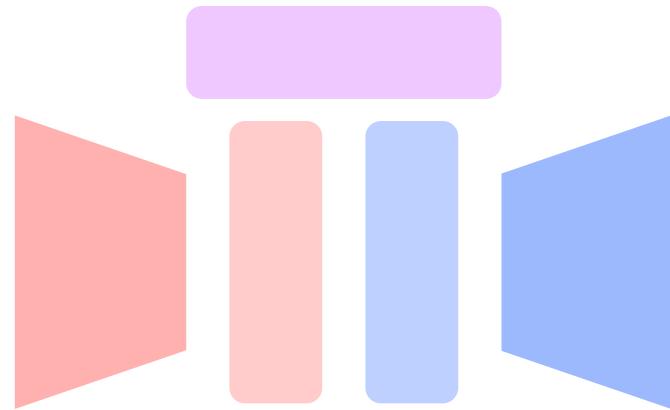
### Sound representation - Spectrogram

- Spectrogram: A stack of spectrums along the time axis

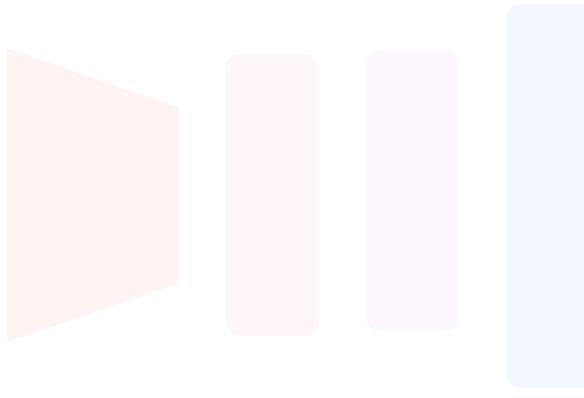


## 3.2 Joint embedding

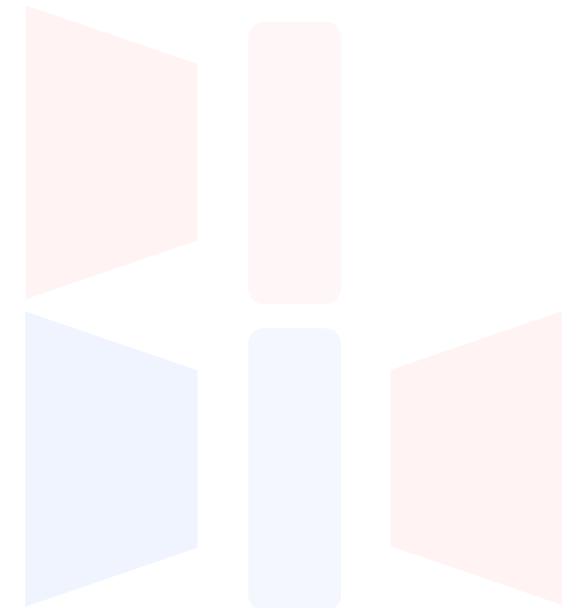
Multi-modal - Audio



Matching



Translating



Referencing

## 3.2 Joint embedding

Multi-modal - Audio

Application - Scene recognition by sound

[Aytar et al., NIPS 2016]



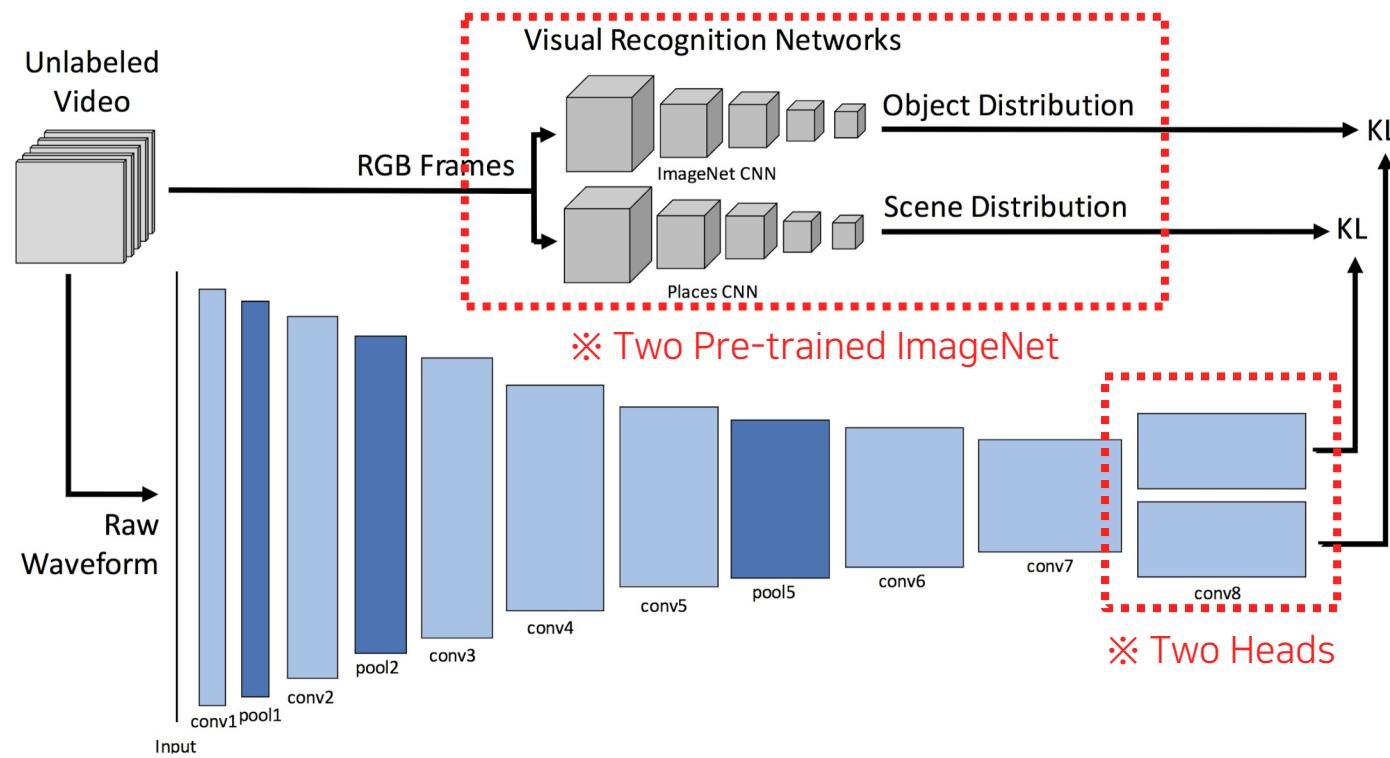
## 3.2 Joint embedding

Multi-modal - Audio

SoundNet

[Aytar et al., NIPS 2016]

- Learn audio representation from synchronized RGB frames in the same videos
- Train by the teacher-student manner
  - Transfer visual knowledge from pre-trained visual recognition models into sound modality



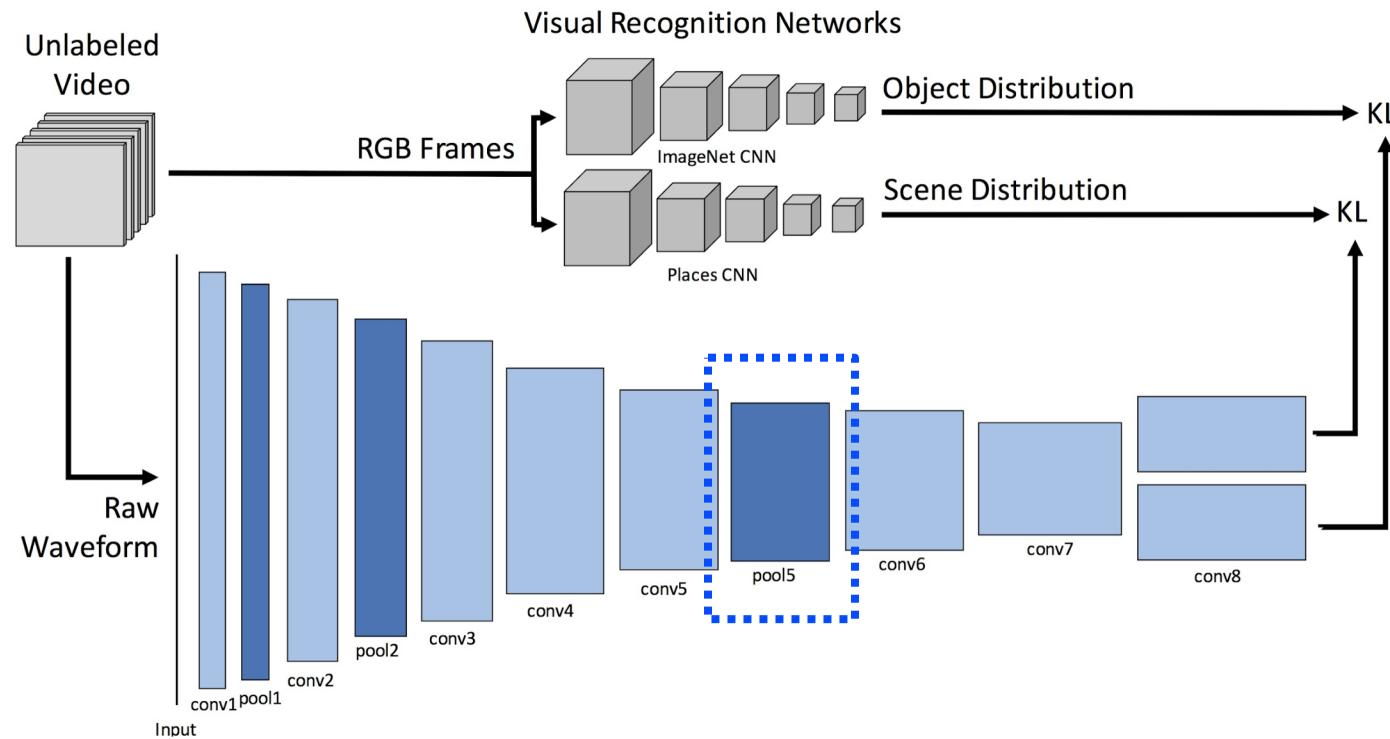
## 3.2 Joint embedding

Multi-modal - Audio

SoundNet

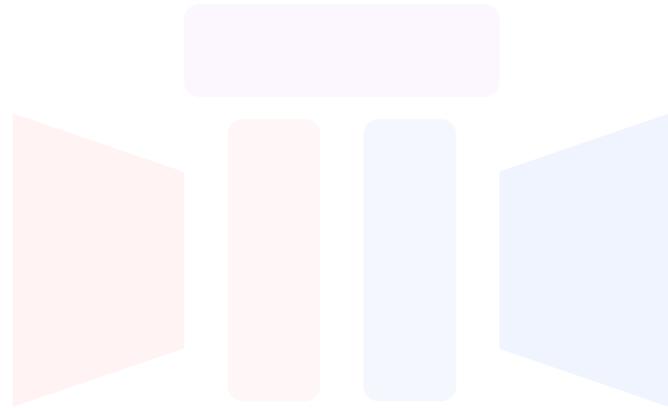
[Aytar et al., NIPS 2016]

- For a target task, the pre-trained internal representation (pool5) is used as features
- Training a classifier with the pool5 feature
  - Instead of the output layer, the pool5 feature posses more generalizable semantic info.

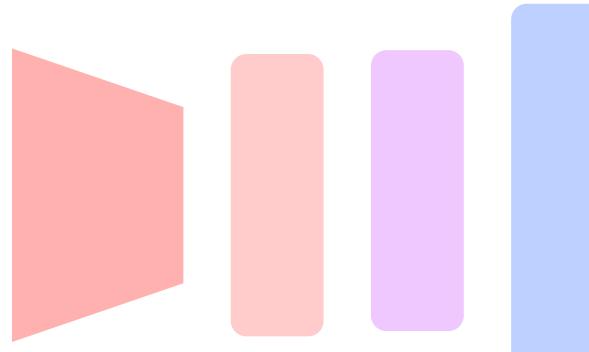


### 3.3 Cross modal translation

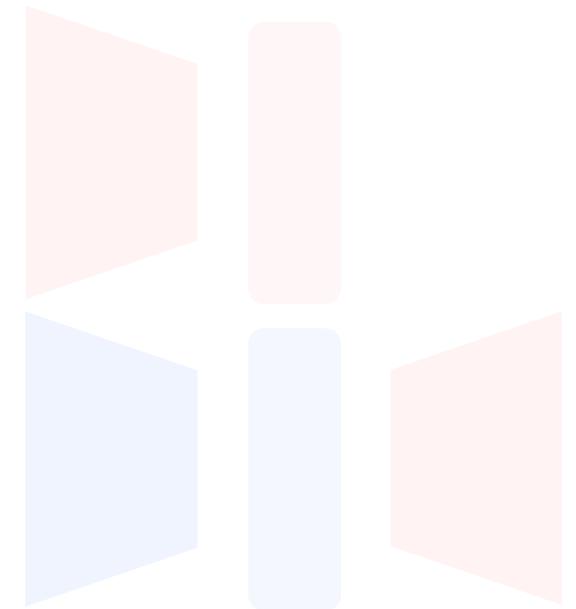
Multi-modal - Audio



Matching



Translating



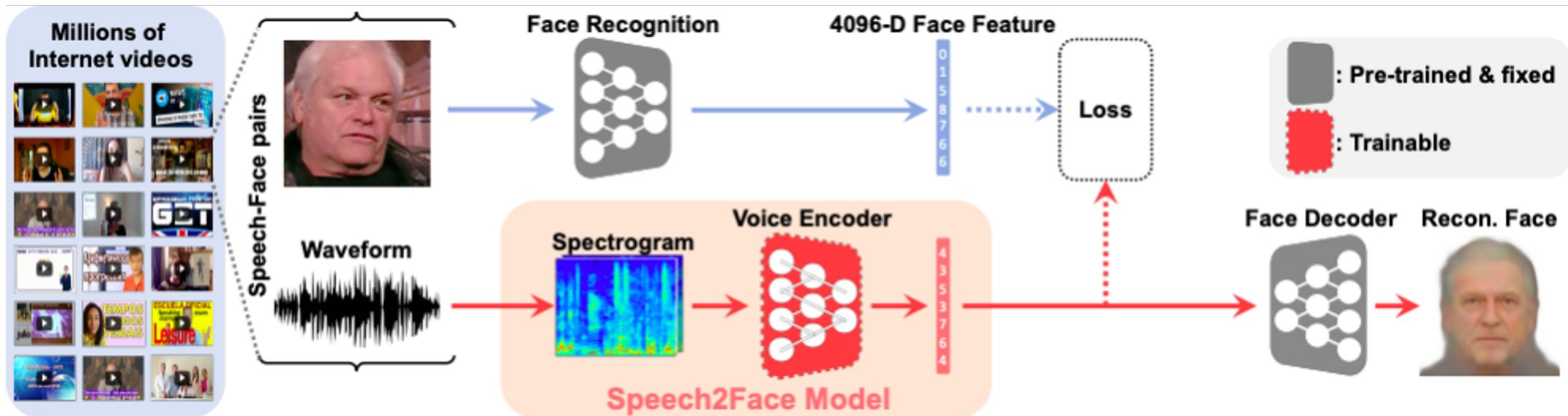
Referencing

### 3.3 Cross modal translation

Multi-modal - Audio

Speech2Face

[Oh et al., CVPR 2019]

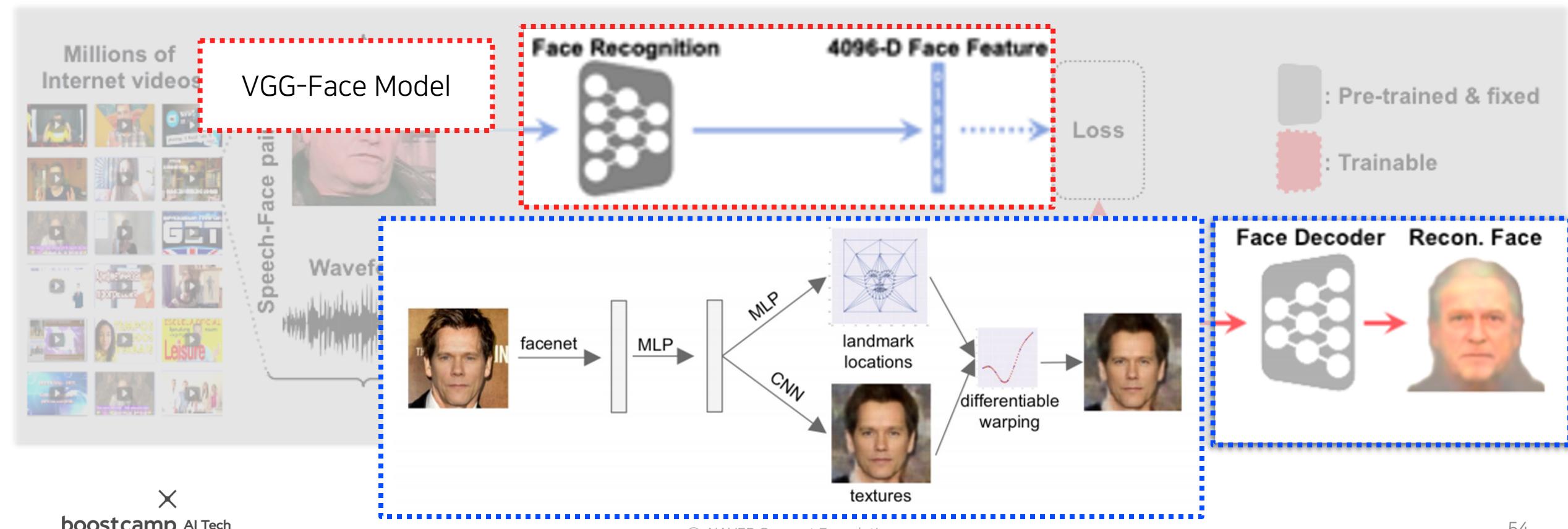


### 3.3 Cross modal translation

Multi-modal - Audio

Speech2Face - Module networks

[Oh et al., CVPR 2019]



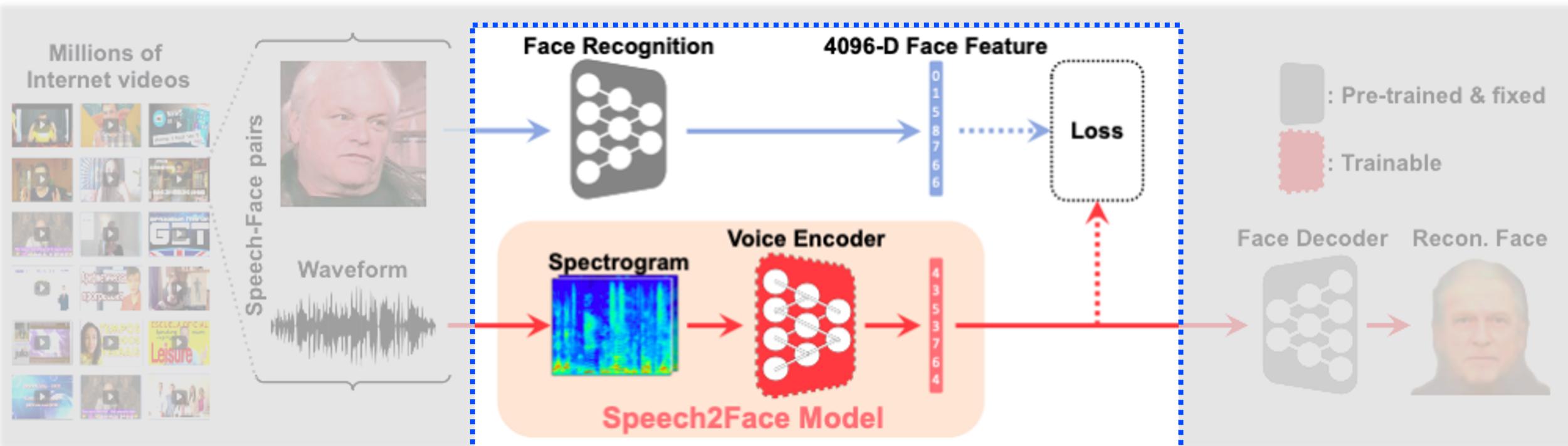
### 3.3 Cross modal translation

Multi-modal - Audio

#### Speech2Face - Training

[Oh et al., CVPR 2019]

- Training by feature matching loss (self-supervised manner) for making features compatible
  - Natural co-occurrence of speaker's speech and facial images

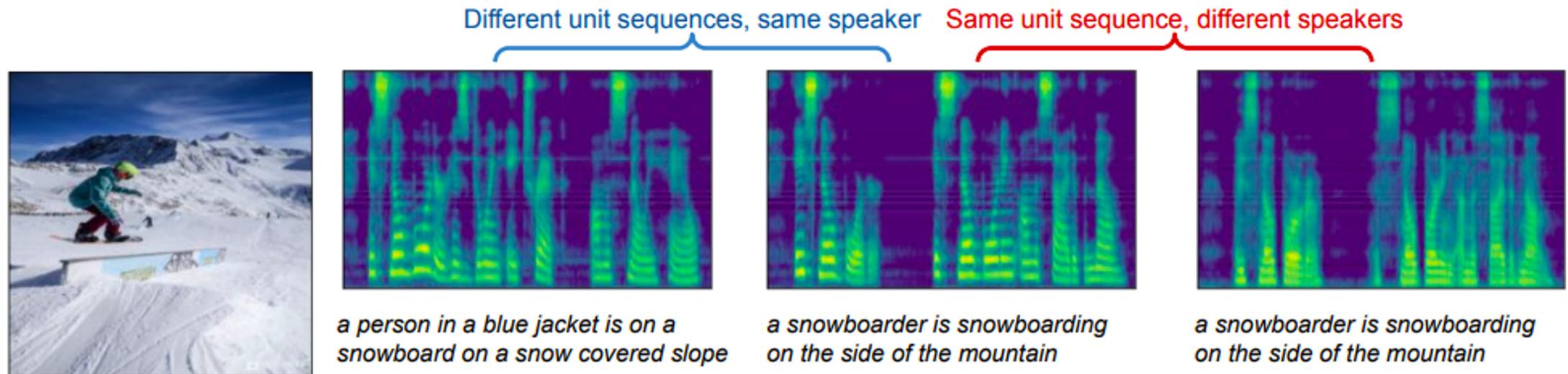


### 3.3 Cross modal translation

Multi-modal - Audio

Application - Image-to-speech synthesis

[Hsu et al., arXiv 2020]



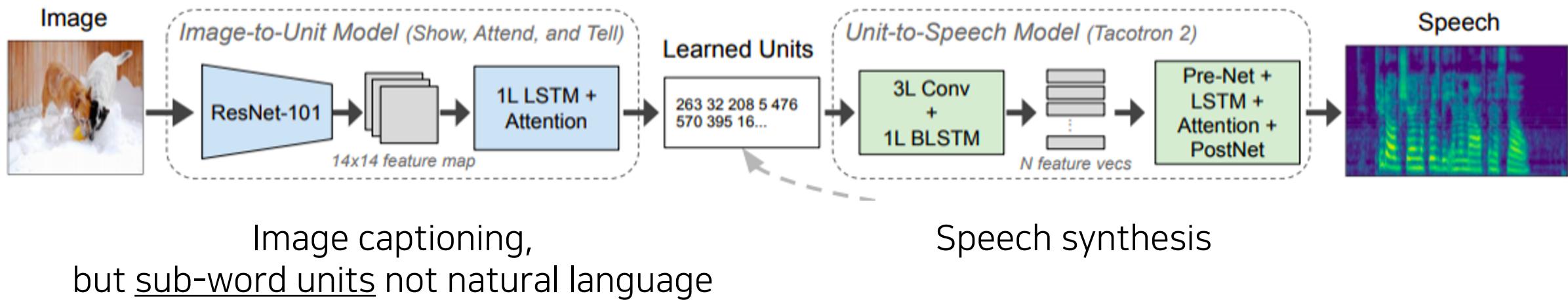
Succeed with discrete unit representation and exhibit robustness to different speakers

### 3.3 Cross modal translation

Multi-modal - Audio

Image-to-speech synthesis – Module networks

[Hsu et al., arXiv 2020]

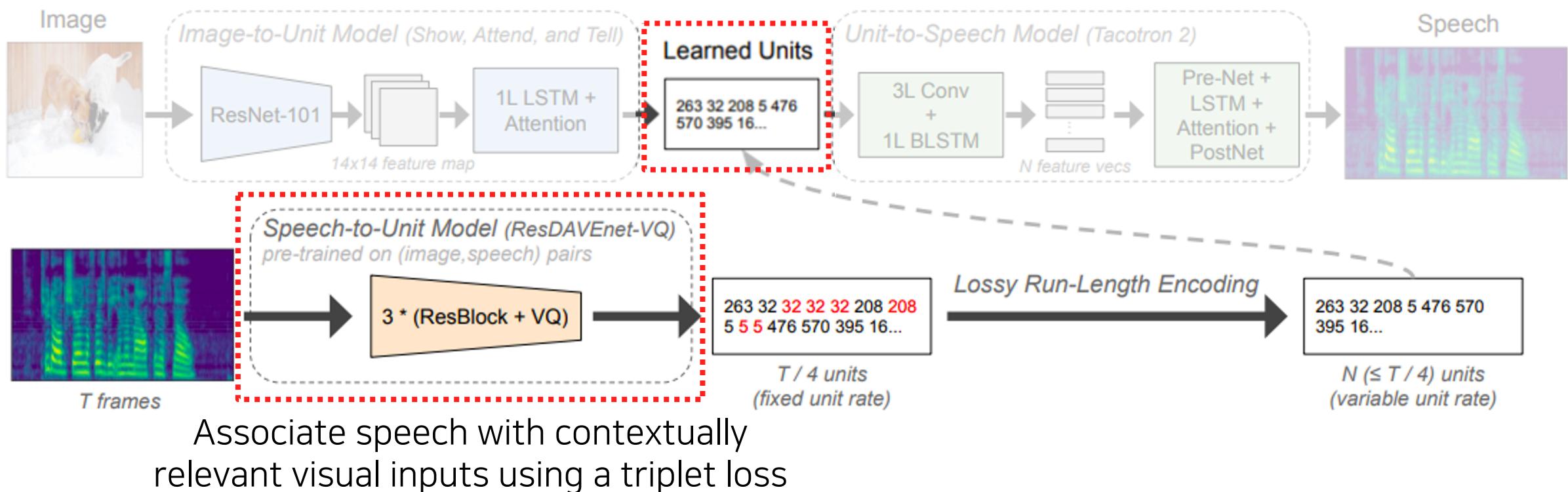


### 3.3 Cross modal translation

Multi-modal - Audio

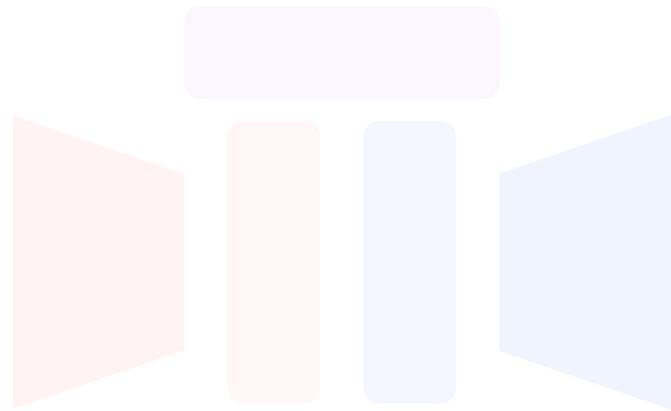
Image-to-speech synthesis – Module networks

[Hsu et al., arXiv 2020]

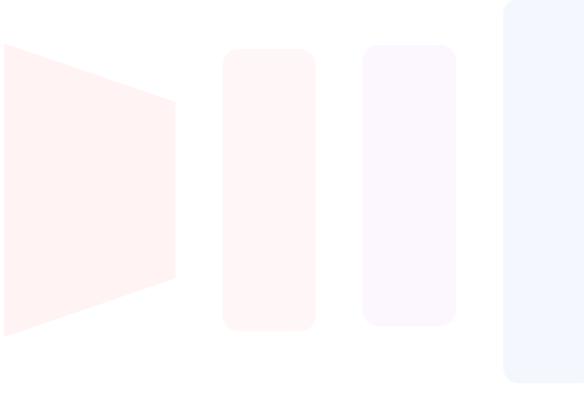


## 3.4 Cross modal reasoning

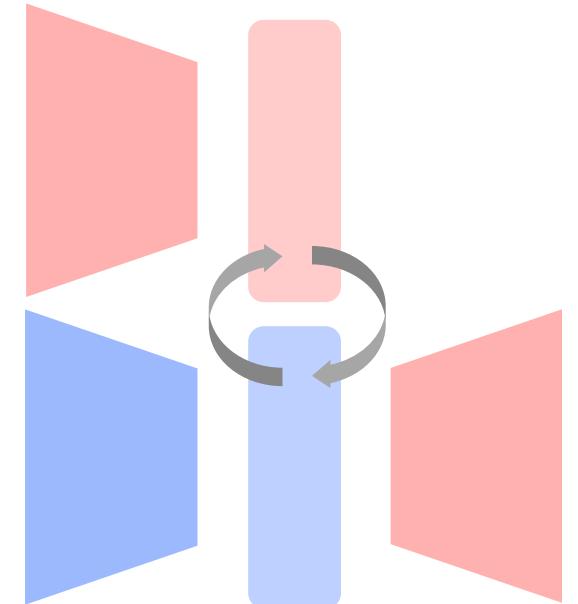
Multi-modal - Audio



Matching



Translating



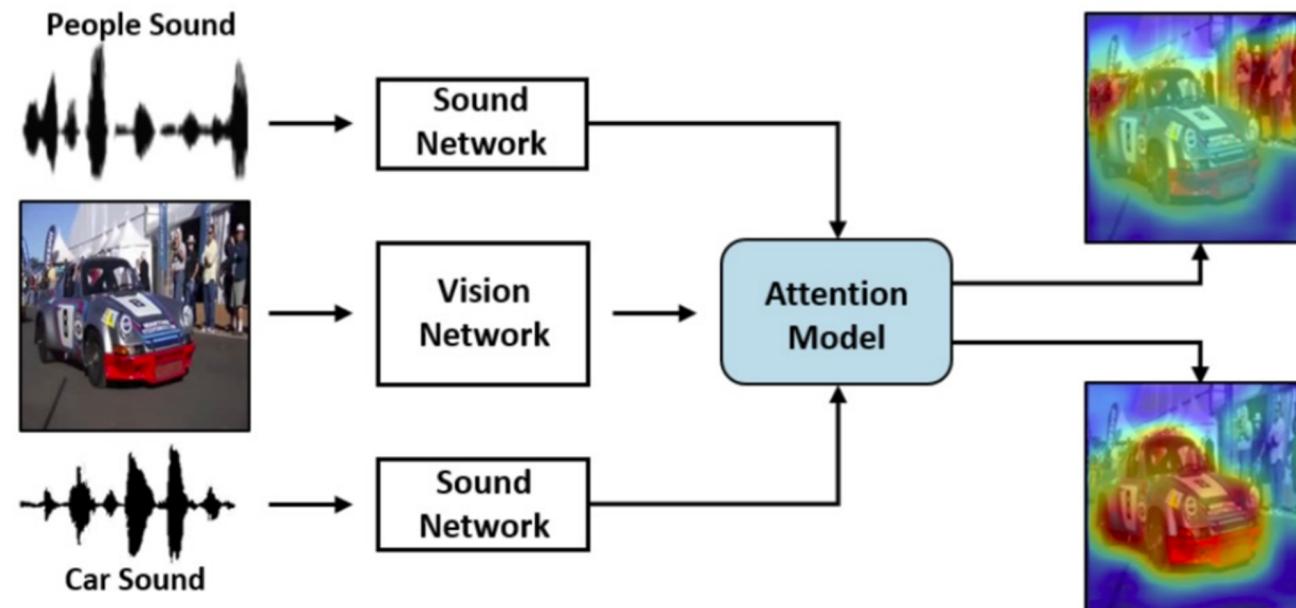
Referencing

## 3.4 Cross modal reasoning

Multi-modal - Audio

Application - Sound source localization

[Senocak et al. CVPR 2018]



## 3.4 Cross modal reasoning

Multi-modal - Audio

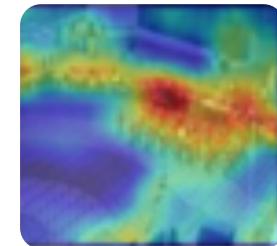
Sound source localization

[Senocak et al. CVPR 2018]

Image



Audio



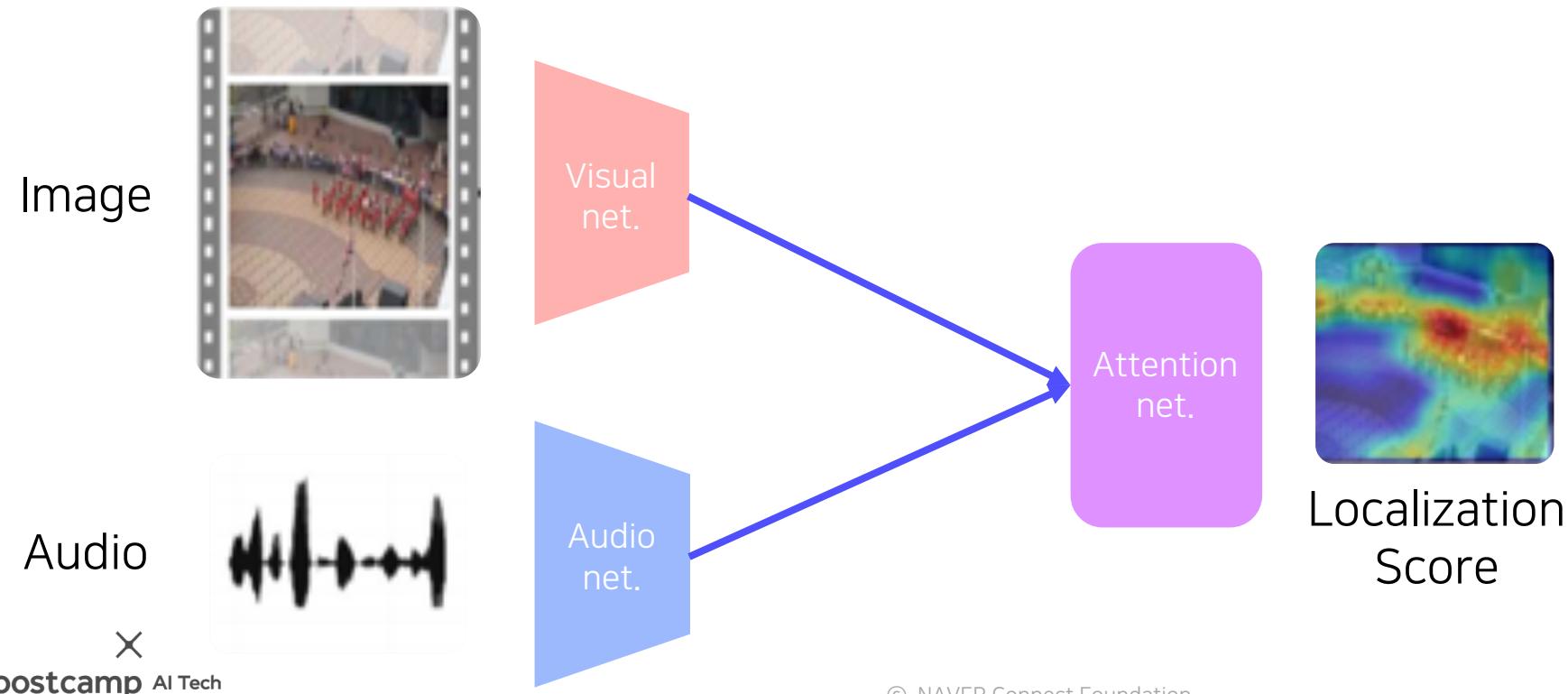
Localization  
Score

## 3.4 Cross modal reasoning

Multi-modal - Audio

Sound source localization

[Senocak et al. CVPR 2018]

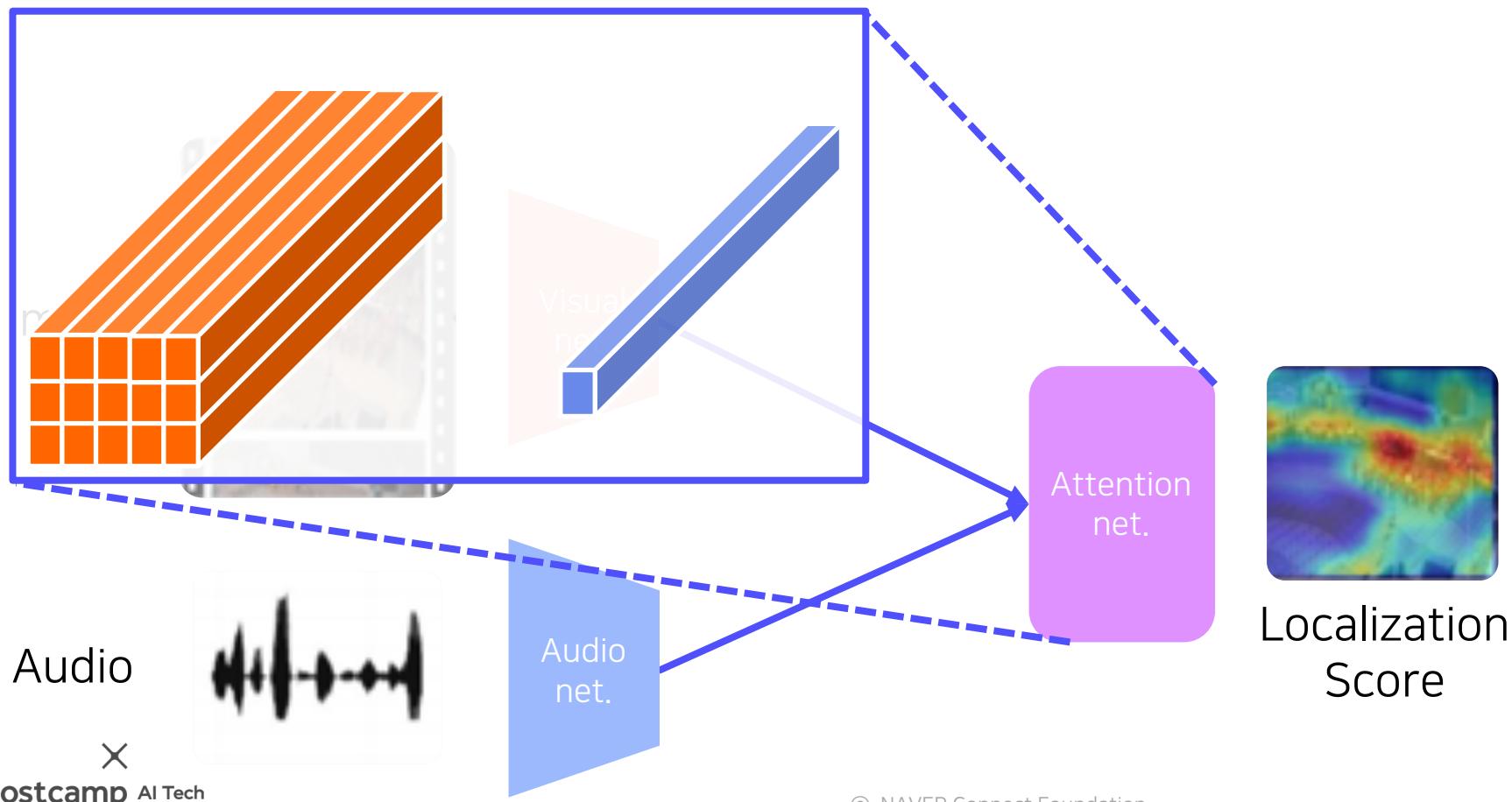


## 3.4 Cross modal reasoning

Multi-modal - Audio

Sound source localization

[Senocak et al. CVPR 2018]

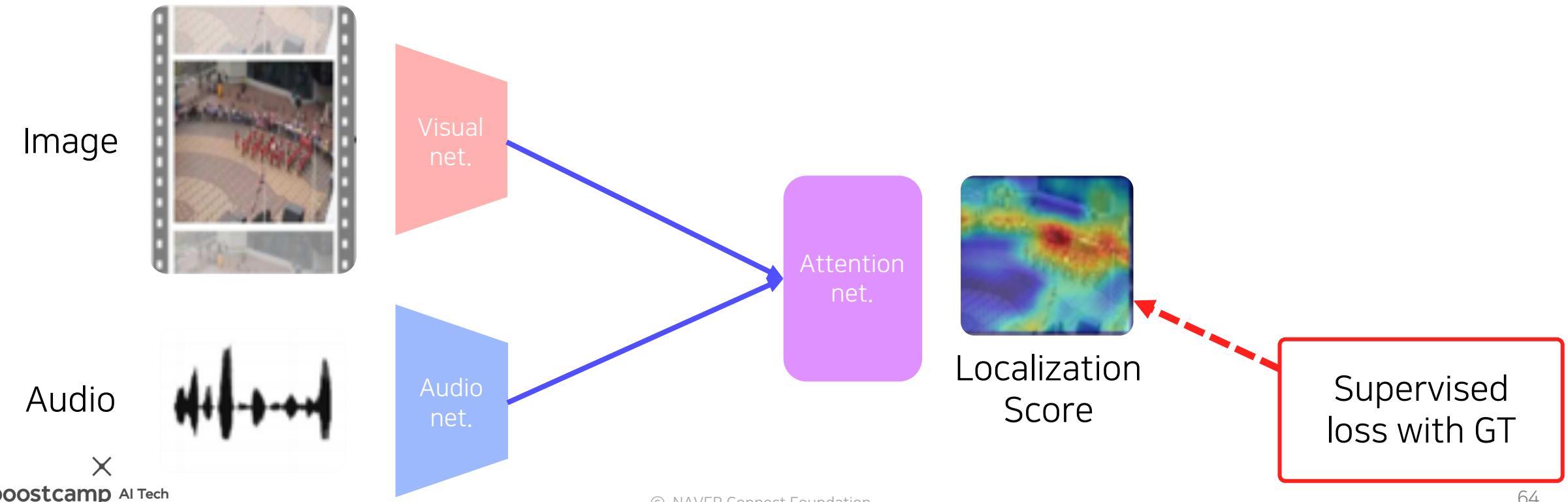


## 3.4 Cross modal reasoning

Multi-modal - Audio

Sound source localization – Fully supervised version

[Senocak et al. CVPR 2018]

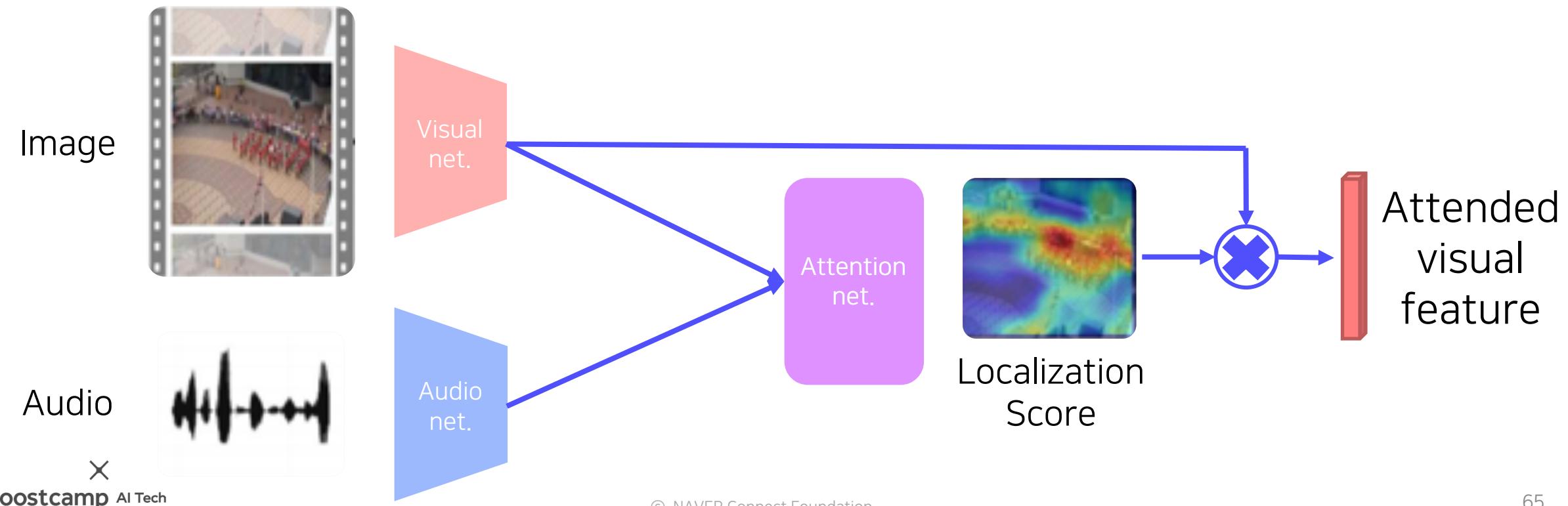


## 3.4 Cross modal reasoning

Multi-modal - Audio

Sound source localization – Unsupervised version

[Senocak et al. CVPR 2018]

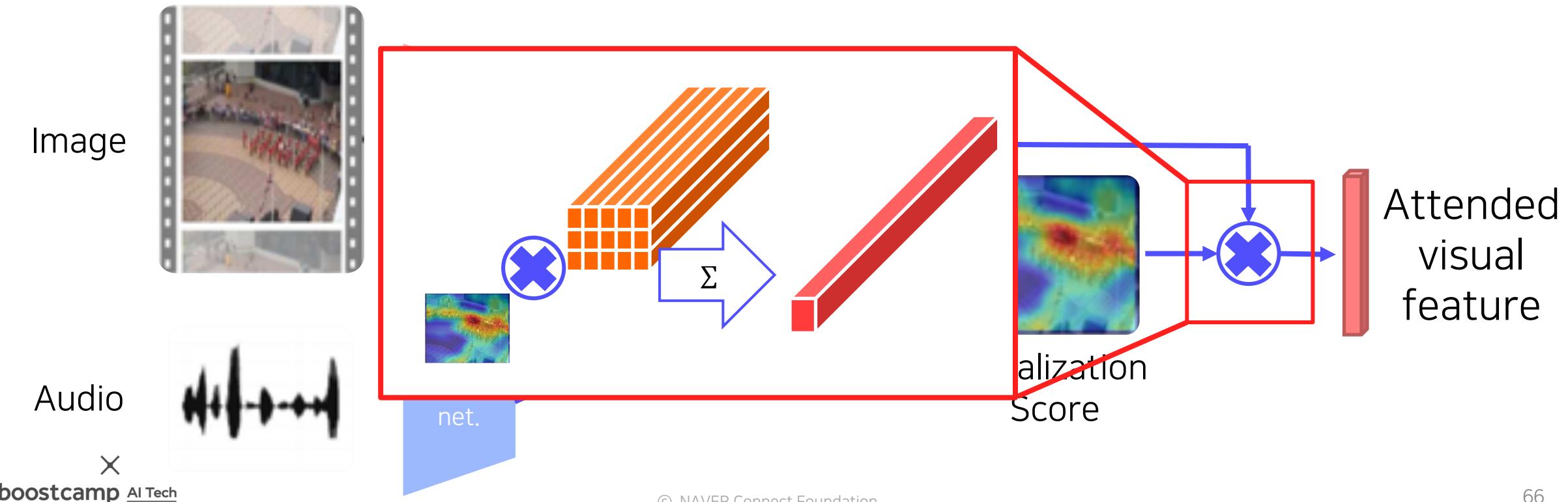


## 3.4 Cross modal reasoning

Multi-modal - Audio

Sound source localization – Unsupervised version

[Senocak et al. CVPR 2018]

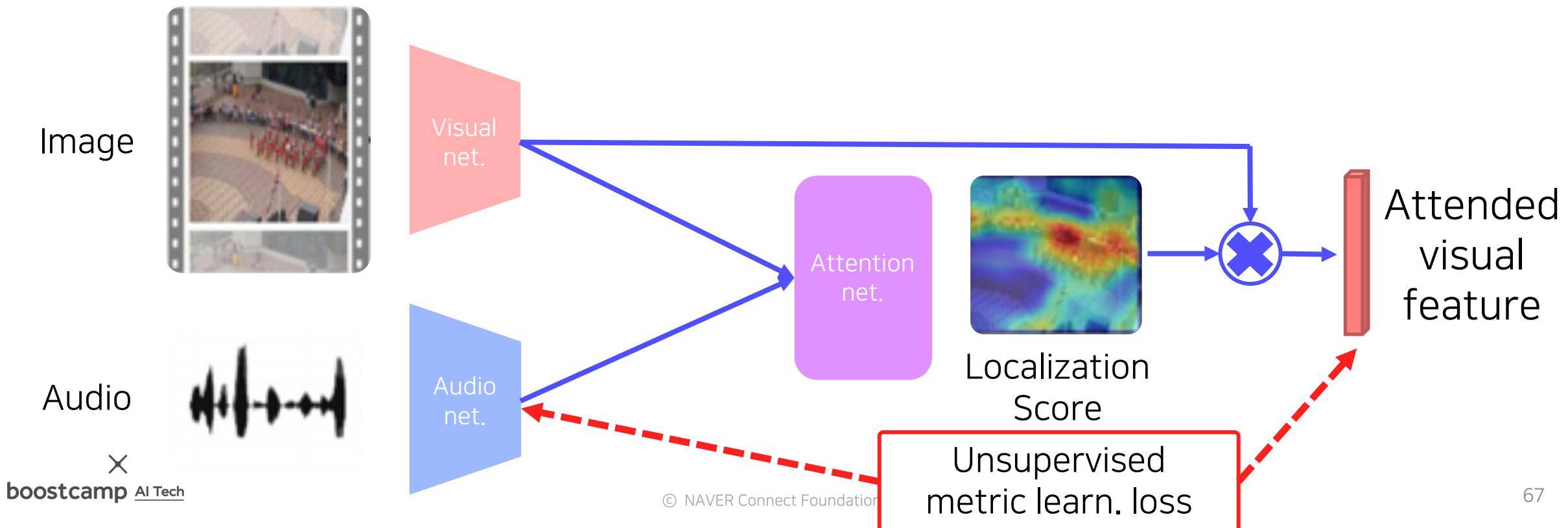


## 3.4 Cross modal reasoning

Multi-modal - Audio

Sound source localization – Unsupervised version

[Senocak et al. CVPR 2018]

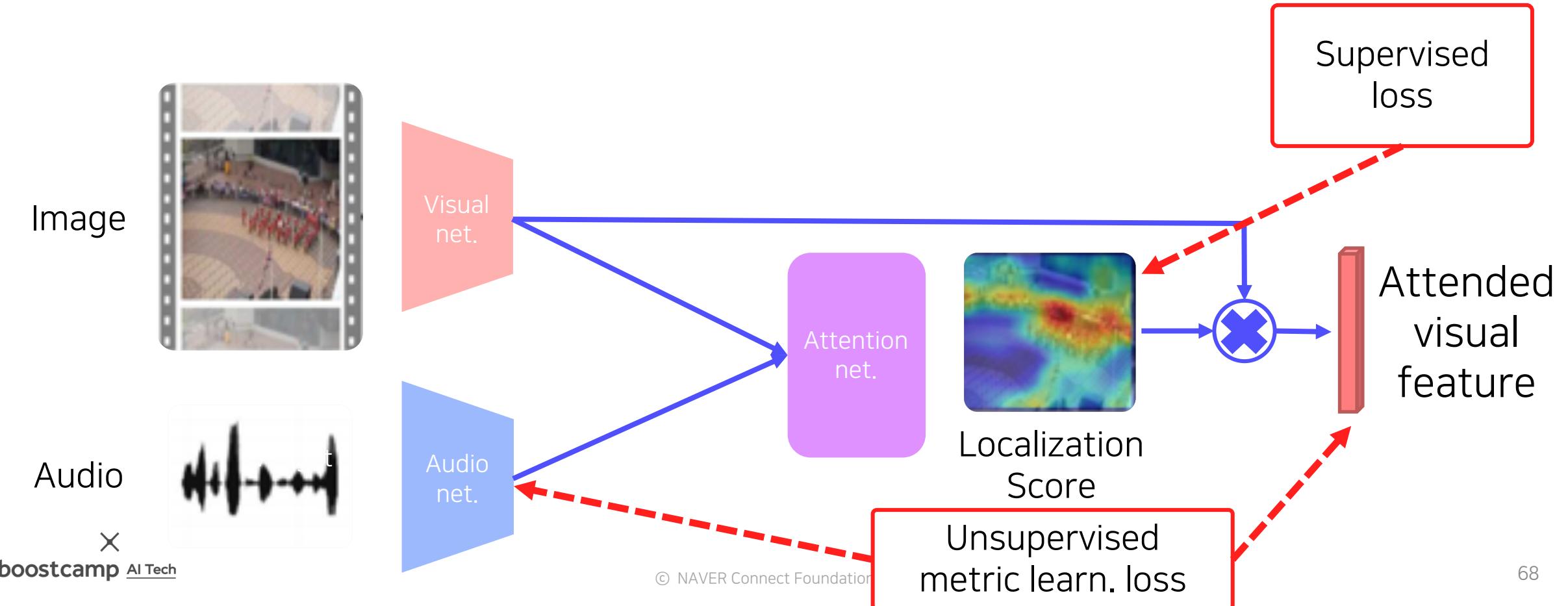


## 3.4 Cross modal reasoning

## Multi-modal - Audio

# Sound source localization – Semi-supervised version

[Senocak et al. CVPR 2018]



## 3.4 Cross modal reasoning

Multi-modal - Audio

Application – Speech separation: Looking to listen at the cocktail party

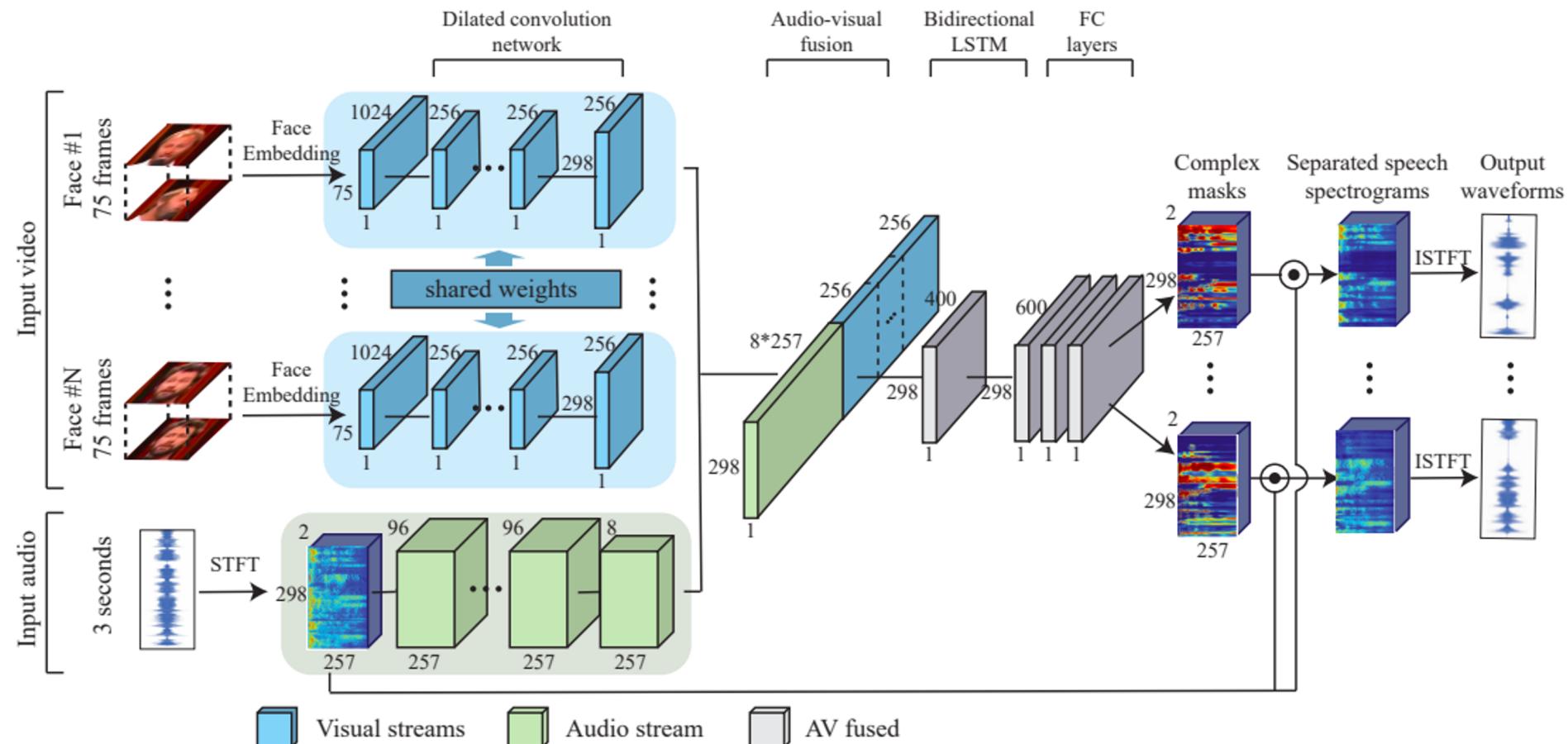
[Ephrat et al., SIGGRAPH 2018]

## 3.4 Cross modal reasoning

Multi-modal - Audio

Looking to listen at the cocktail party

[Ephrat et al., SIGGRAPH 2018]

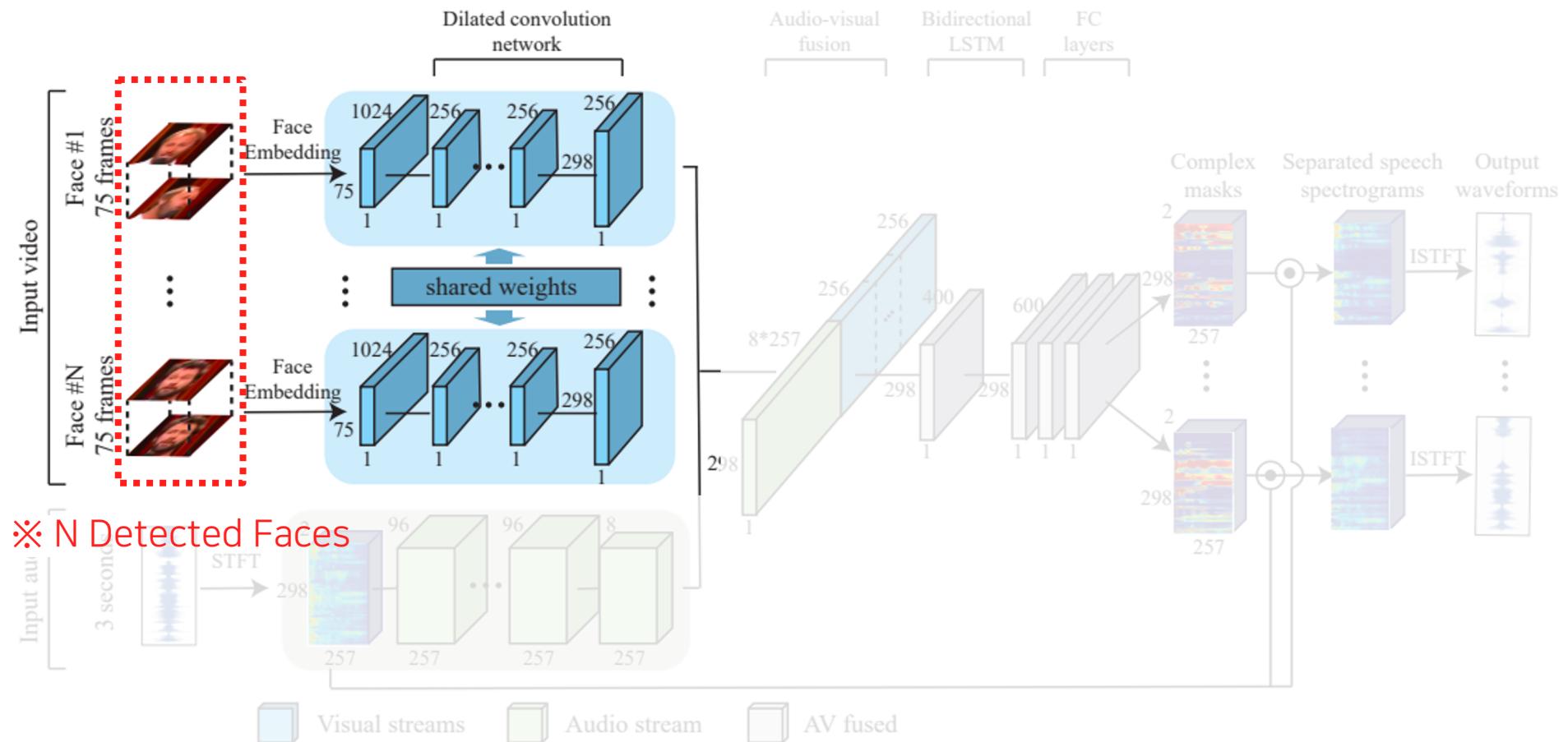


## 3.4 Cross modal reasoning

Multi-modal - Audio

Looking to listen at the cocktail party - Visual stream

[Ephrat et al., SIGGRAPH 2018]

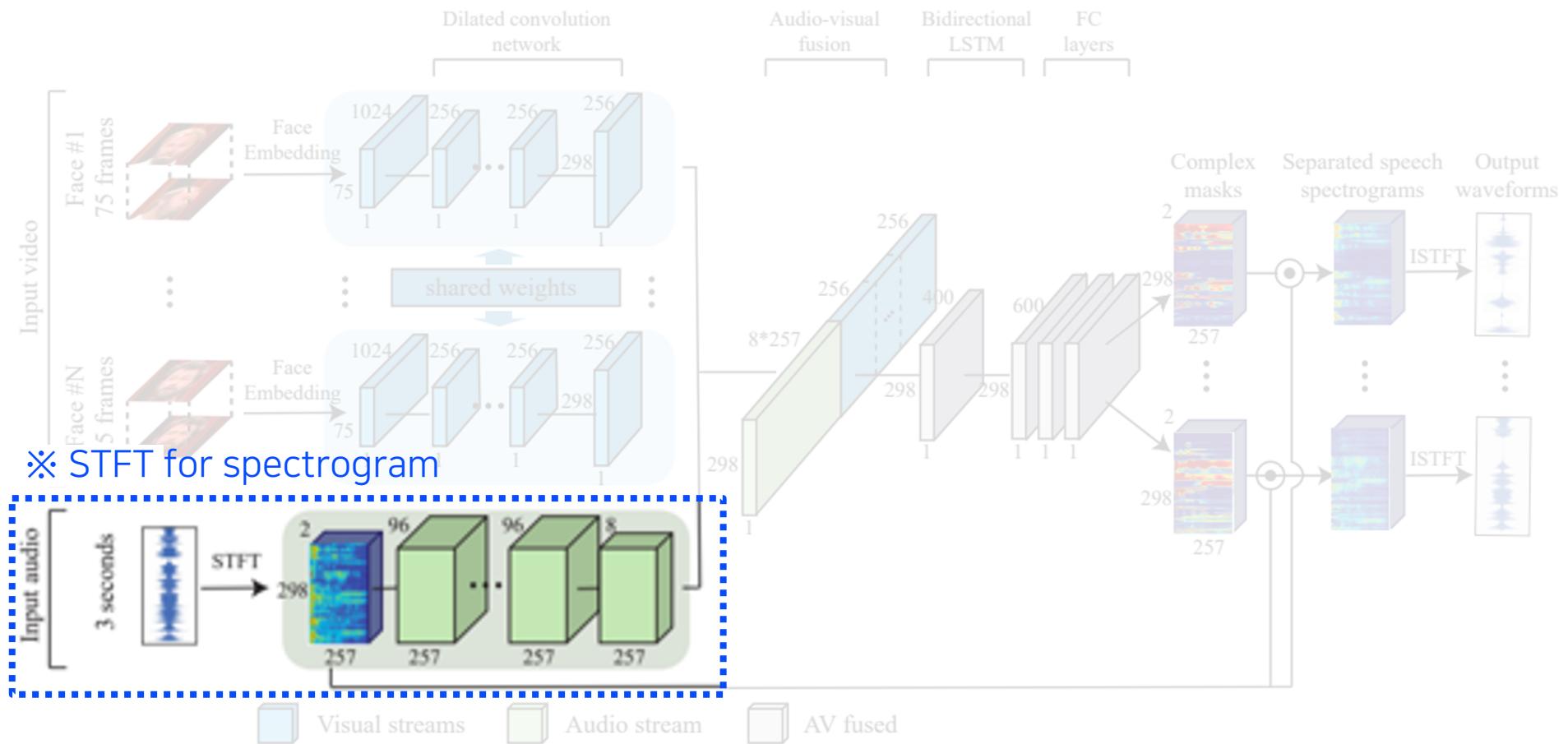


### 3.4 Cross modal reasoning

Multi-modal - Audio

Looking to listen at the cocktail party - Audio stream

[Ephrat et al., SIGGRAPH 2018]



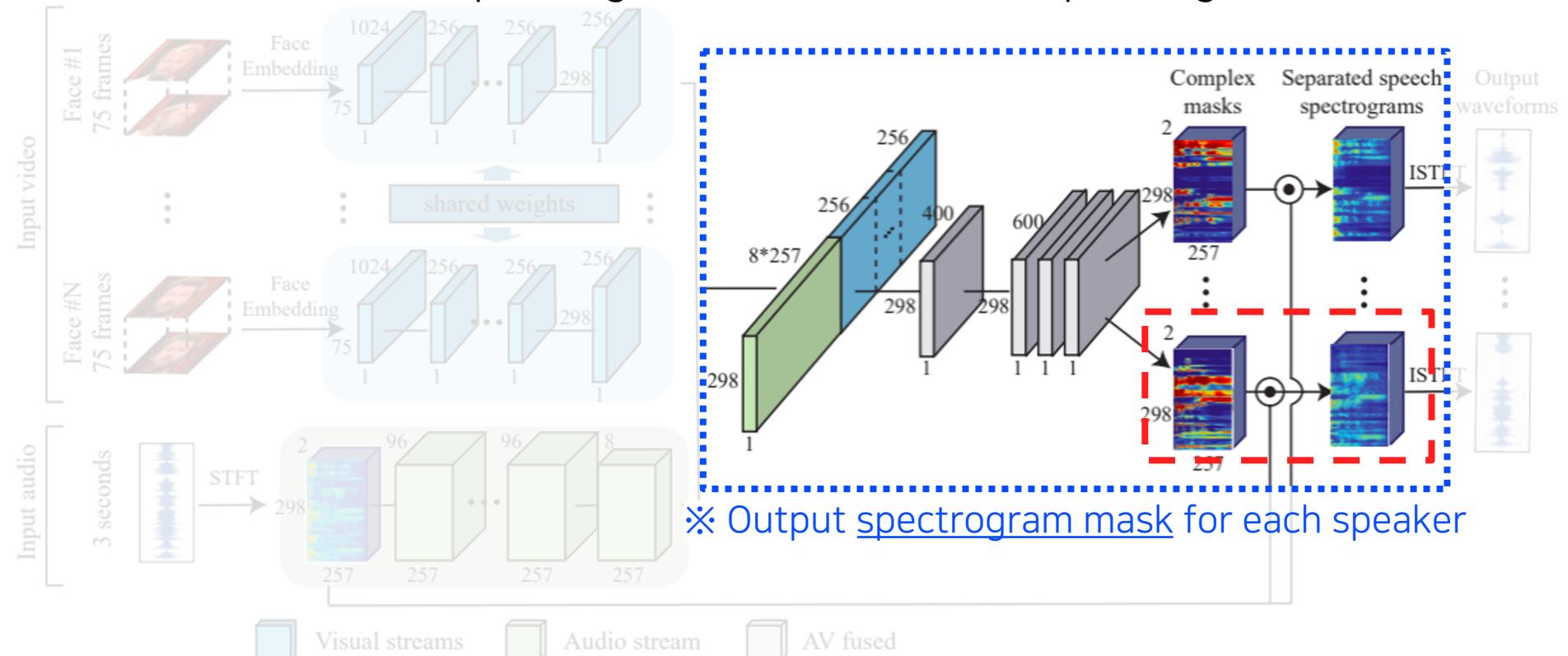
## 3.4 Cross modal reasoning

Multi-modal - Audio

Looking to listen at the cocktail party - Audio-visual fusion

[Ephrat et al., SIGGRAPH 2018]

- Training data: synthetically generated by combining two clean speech videos
- Loss: L2 loss between “clean spectrogram” and “enhanced spectrogram”



## 3.4 Cross modal reasoning

Multi-modal - Audio

Lip movements generation - Synthesizing Obama example

[Suwajanakorn et al., SIGGRAPH 2017]

Speech source

Synthesized Obama

# Conclusion. Beyond image, text and audio

Multi-modal - Audio

---

Autopilot - Tesla self-driving

---

## 1. Multi-modal learning overview

- Wang et al., What Makes Training Multi-Modal Classification networks Hard?, CVPR 2020

## 2. Multi-modal (1) – Text

- Srivastava and Slakhutdinov, Multimodal Learning with Deep Boltzmann Machines, JMLR 2014
- Marin et al., Recipe 1M+: A Dataset for Learning Cross-Modal Embeddings for Cooking Recipes and Food Images, TPAMI 2019
- Vinyals et al. , Show and Tell: A Neural Image Caption Generator, CVPR 2015
- Xu et al., Show, attend and tell: Neural image caption generation with visual attention, ICML 2015
- Reed et al., Generative Adversarial Text to Image Synthesis, ICML 2016
- Antol et al., VQA: Visual Question Answering, ICCV 2015

# Reference

---

## 3. Multi-modal (2) - Audio

- Aytar et al., SoundNet: Learning Sound Representations from Unlabeled Video, NIPS 2016
- Senocak et al., Learning to Localize Sound Sources in Visual Scenes: Analysis and Applications, CVPR 2018
- Oh et al., Speech2Face: Learning the Face Behind a Voice, CVPR 2019
- Hsu et al., Text-Free Image-to-Speech Synthesis Using Learned Segmental Units, arXiv 2020
- Ephrat et al., Looking to Listen at the Cocktail Party: A Speaker-Independent Audio-Visual Model for Speech Separation, SIGGRAPH 2018
- Suwajanakorn et al., Synthesizing Obama: Learning Lip Sync from Audio, SIGGRAPH 2017
- Chen et al., Lip Movements Generation at a Glance, ECCV 2018

# End of Document

## Thank You.

상위 카테고리 입력란