# Multimodal Machine Learning

## Lecture 10.1: Fusion, co-learning and new trends

**Louis-Philippe Morency**

* Original version co-developed with Tadas Baltrusaitis

# Administrative Stuff

# Piazza Live Q&A



Please share your questions and comments on Piazza Live Q&A

➡ Live responses by your TAs and follow-up by the instructor after the main lecture

Language Technologies Institute

Carnegie Mellon University

# Lecture Schedule

| Classes | Tuesday Lectures | Thursday Lectures |
|---|---|---|
| **Week 7**<br>10/13 & 10/15 | **Alignment and translation**<br>• Neural Module networks<br>• Connectionist temporal classification | **Probabilistic graphical models**<br>• Dynamic Bayesian networks<br>• Coupled and factor HMMs |
| **Week 8**<br>10/20 & 10/22 | **Discriminative graphical models**<br>• Conditional random fields<br>• Continuous and fully-connected CRFs | **Neural Generative Models**<br>• Variational auto-encoder<br>• Generative adversarial networks |
| **Week 9**<br>10/27 & 10/29 | **Reinforcement learning**<br>• Markov decision process<br>• Q learning and policy gradients | **Multimodal RL**<br>• Deep Q learning<br>• Multimodal applications |
| **Week 10**<br>11/3 & 11/5 | **Fusion and co-learning**<br>• Multi-kernel learning and fusion<br>• Few shot learning and co-learning | **New research directions**<br>• Recent approaches in multimodal ML |
| **Week 11**<br>11/10 & 11/12 | ***Mid-term project assignment** (live working sessions instead of lectures)* | |

**Midterm project assignment**
Presentations due Friday 11/13
Reports due Sunday 11/15
Peer feedback due Sunday 11/22

Language Technologies Institute

Carnegie Mellon University

# Lecture Schedule

| Classes | Tuesday Lectures | Thursday Lectures |
|---------|------------------|-------------------|
| **Week 12**<br>11/17 & 11/19 | **Embodied Language Grounding**<br>• Connecting Language to Action<br>• Guest lecture: Yonatan Bisk | **Multimodal language acquisition**<br>• Learning from multimodal data<br>• Guest lecture: Graham Neubig |
| **Week 13**<br>11/24 & 11/26 | ***Thanksgiving week*** *(no lectures)* | |
| **Week 14**<br>12/1 & 12/3 | **Learning to connect text and images**<br>• Discourse approaches, text & images<br>• Guest lecture: Malihe Alikhani | **Bias and fairness**<br>• Computational ethics<br>• Guest lecture: Yulia Tsvetkov |
| **Week 15**<br>12/8 & 12/10 | ***Final project assignment*** *(live working sessions instead of lectures)* | |

**Final project assignment**
Presentations due Friday 12/11
Reports due Sunday 12/13

Language Technologies Institute

Carnegie Mellon University

# Next Week Schedule

**Tuesday (11/10) 3pm-6pm:** Live office hours with LP

- Signup on Calendly for meeting timeslot (see next slide)
- Use the same Zoom link (waiting room will be activated)

**Thursday (11/12) :** No lecture

**Friday (11/13) 8pm:** deadline for presentations

- Submit on Gradescope (slides) and Box (video)

**Sunday (11/15) 8pm:** deadline for reports

- Submit on Gradescope

**Sunday (10/9) 8pm:** Deadline for student feedback

No reading assignment for Week 11

Reading assignment for Week 12 (starting Monday 11/16)

# Signup Sheet for LP's Office Hours

**Tuesday (11/10) 3pm-6pm**

Sign-up using Calendly:

https://calendly.com/morency/student-meetings

- One meeting per team
- Each meeting 10mins (-ish)
- Same Zoom link as lectures
    - Waiting room will be activated

**Language Technologies Institute**

**Carnegie Mellon University**

# Multimodal Machine Learning

## Lecture 10.1: Fusion, co-learning and new trends

**Louis-Philippe Morency**

* Original version co-developed with Tadas Baltrusaitis

# Lecture Objectives

- Quick recap: multimodal fusion

- Model-agnostic fusion
    - Multimodal fusion architecture search

- Fusion and kernel function
    - Transformers through the lens for kernel
    - Multiple Kernel Learning

- Co-learning
    - Paired and weakly-paired data

- Research trends in Multimodal ML *NEW* **papers**
    - Few-shot and weakly supervised learning
    - Multi-lingual multimodal grounding

# Quick Recap: Multimodal Fusion

Language Technologies Institute

Carnegie Mellon University

# Multimodal fusion

- Process of joining information from two or more modalities to perform a prediction

- Examples
  - Audio-visual speech recognition
  - Audio-visual emotion recognition
  - Multimodal biometrics
  - Speaker identification and diarization
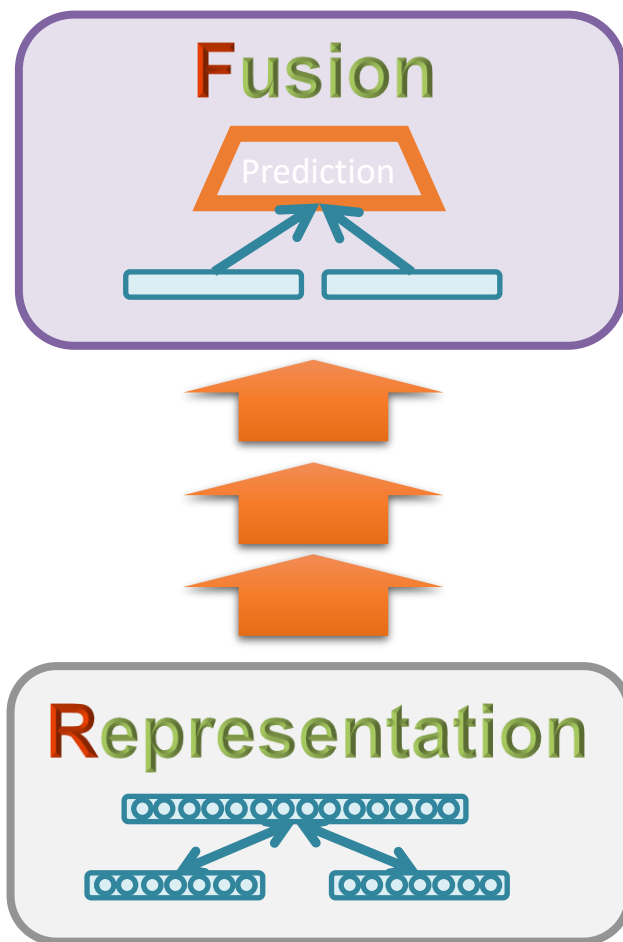  - Visual/Media Question answering

(a) answer-phone      (a) get-out-car      (a) fight-person

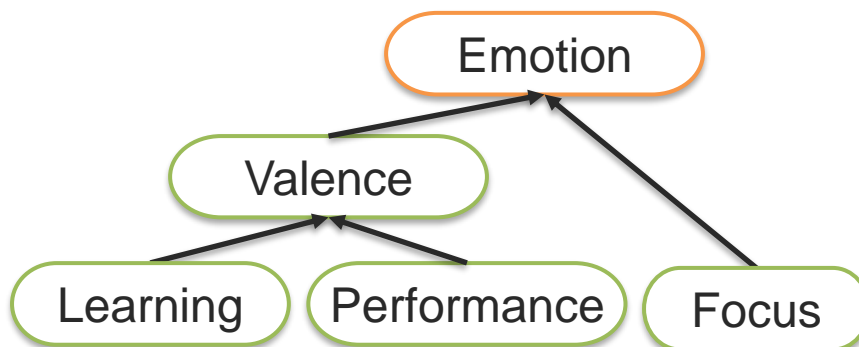# Fusion – Probabilistic Graphical Models

**Fusion**

Prediction

← **Domain knowledge**

a) Latent sub-structure

Emotion

Valence

Learning    Performance    Focus

**Representation**

b) Structured output prediction

$Emotion_{t-1}$    $Emotion_t$

# Graphical Model: Learning Multimodal Structure

Modality-*private* structure

- Internal grouping of observations

Modality-*shared* structure

- Interaction and synchrony

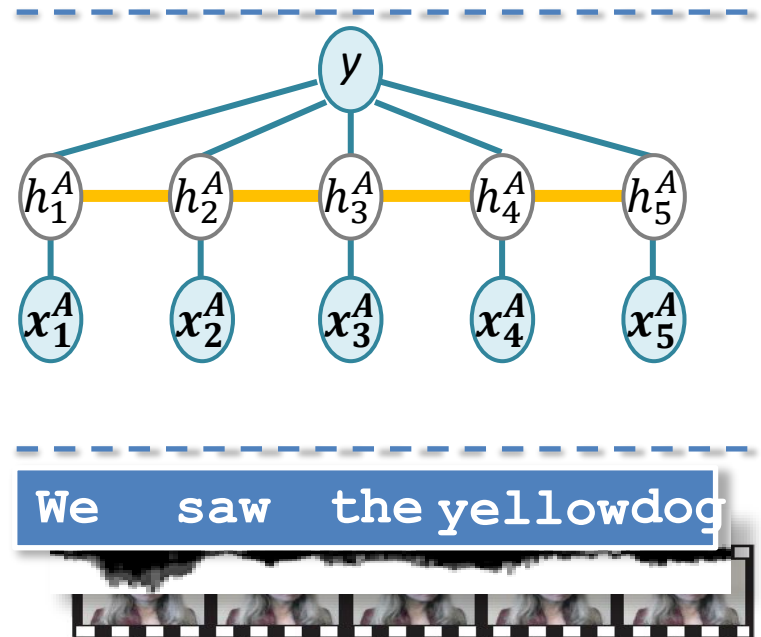Language Technologies Institute
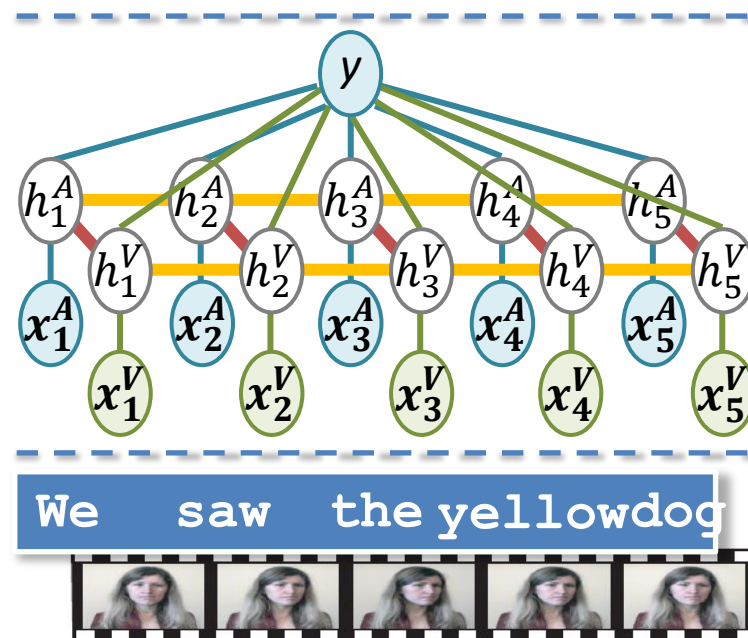
Carnegie Mellon University

# Graphical Model: Learning Multimodal Structure

Modality-*private* structure

  • Internal grouping of observations

Modality-*shared* structure

  • Interaction and synchrony



$$p(y \mid \boldsymbol{x}^A, \boldsymbol{x}^V; \boldsymbol{\theta}) = \sum_{\boldsymbol{h}^A, \boldsymbol{h}^V} p(y, \boldsymbol{h}^A, \boldsymbol{h}^V \mid \boldsymbol{x}^A, \boldsymbol{x}^V; \boldsymbol{\theta})$$

➤ Approximate inference using loopy-belief

Language Technologies Institute
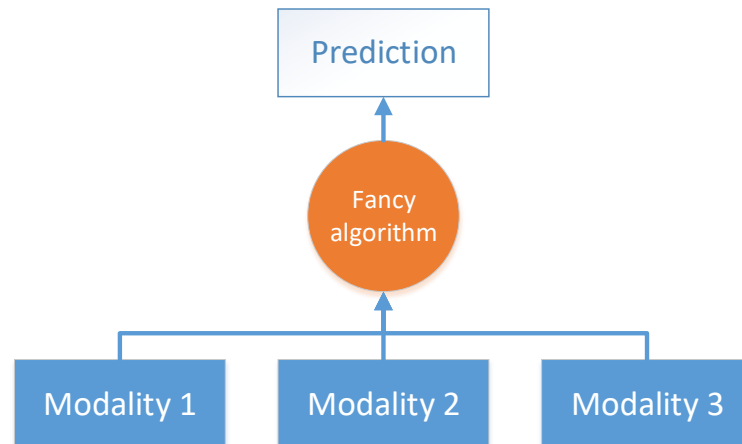
Carnegie Mellon University

# Multimodal Fusion

"Model-agnostic" fusion:

- Early and late fusion
- Fusion architecture search

Intermediate fusion (aka model-based):

- Neural Networks
- Graphical models
- Kernel Methods

# Model-free Fusion

# Model-agnostic approaches – early fusion
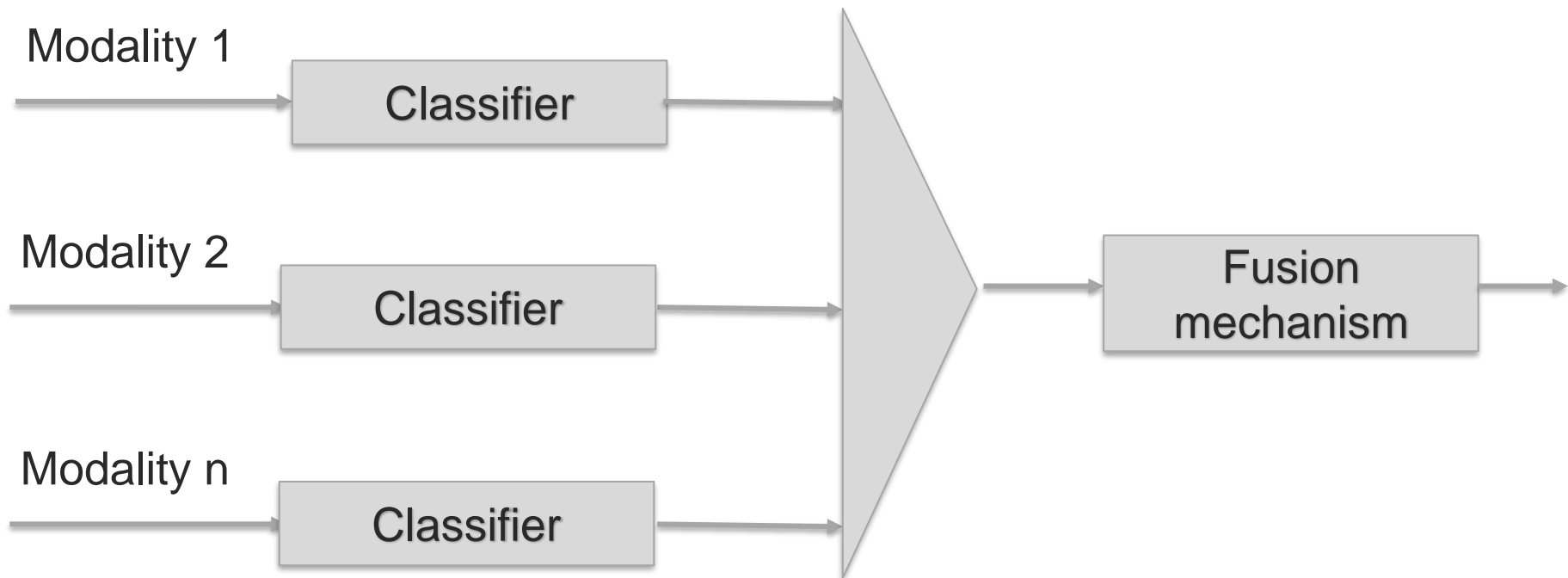
Modality 1

Modality 2

Modality n

Classifier

- Easy to implement – just concatenate the features
- Exploit dependencies between features
- Can end up very high dimensional
- More difficult to use if features have different granularities

# Model-agnostic approaches – late fusion

Modality 1

Classifier

Modality 2

Classifier

Modality n

Classifier

Fusion mechanism

- Train a unimodal predictor and a multimodal fusion one
- Requires multiple training stages
- Do not model low level interactions between modalities
- Fus~~ion~~ ~~a~~pproach

What should be the Fusion Mechanism for multi-layer neural classifiers?

# Late Fusion on Multi-Layer Unimodal Classifiers

Unimodal classifier 1



Unimodal classifier 2

## What layer(s) should we fuse?

One of the last layers?



$L=1$

Or more than one layer?



$L=2$

Trying all combinations may be computationally expensive…

# Multimodal Fusion Architecture Search (MFAS)
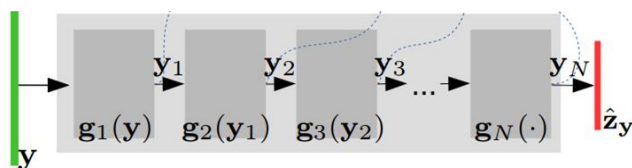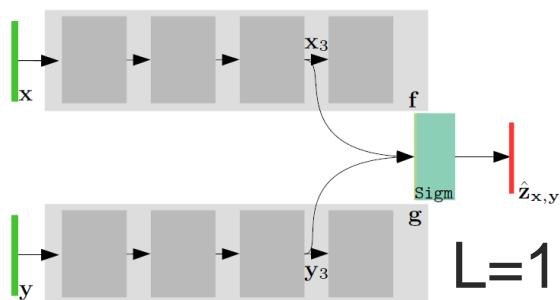
**NEW**-ish **paper**

**Proposed solution:** Explore the search space with *Sequential Model-Based Optimization*

➡ Start with simpler models first (all L=1 models) and iteratively increase the complexity (L=2, L=3,…)

➡ Use a *surrogate* function to predict performance of unseen architectures

➡ e.g., the performance of all the L=1 models should give us an idea of how well the L=2 models will perform

"Perez-Rua, Vielzeuf, Pateux, Baccouche, Frederic Jurie,MFAS: Multimodal Fusion Architecture Search, CVPR 2019

Language Technologies Institute

Carnegie Mellon University

# Multimodal Fusion Architecture Search (MFAS)

**Basic building block:** a "fusion layer" unit

Unimodal classifier 1



Fusion layer unit

Unimodal classifier 2

With three hyper-parameters:
a) Layer index for modality 1
b) Layer index for modality 2
c) Activation/fusion function

"Perez-Rua, Vielzeuf, Pateux, Baccouche, Frederic Jurie,MFAS: Multimodal Fusion Architecture Search, CVPR 2019

# Multimodal Fusion Architecture Search (MFAS)

**Dataset:** Audio-Visual MNIST

Example of discovered fusion architecture with MFAS:

LeNet trained on
spoken digits

LeNet trained on
visual digits



"Perez-Rua, Vielz

What should be the Fusion Mechanism for
variable length unimodal classifier?

Languag                                                                                    versity

# Memory-Based Fusion



Local evidences

Memory

$t-1$    $t$    $t+1$    $t+2$

➢ This model can also be trained end-to-end.

[Zadeh et al., Memory Fusion Network for Multi-view Sequential Learning, AAAI 2018]

Language Technologies Institute

Carnegie Mellon University

# Local Fusion and Kernel Functions

Language Technologies Institute

Carnegie Mellon University

# Recap: Transformer Self-Attention

# Transformer Self-Attention

Language Technologies Institute

Carnegie Mellon University

# Transformer's Attention Function

$h_1$

Scale dot-product attention

$\Sigma$

$\alpha_{2,1}$
$\alpha_{1,1}$

$v_1$

$q_1$  $k_1$

$W_q$  $W_k$  $W_v$

$x_1$

**Scale dot-product attention:**

$$\boldsymbol{\alpha} = softmax\left(\frac{\boldsymbol{x_q W_q (x_k W_k)}^T}{\sqrt{d_k}}\right)$$

This attention function is a similarity function. This is related to kernel function…

Language Technologies Institute

Carnegie Mellon University

# What is a Kernel function?

**A kernel function:** Acts as a similarity metric between data points

$$K(x_i, x_j) = \phi(x_i)^T \phi(x_j) = \langle \phi(x_i), \phi(x_j) \rangle, \text{ where } \phi: D \rightarrow Z$$

- Kernel function performs an inner product in feature map space $\phi$
- Inner product (a generalization of the dot product) is often denoted as $\langle .,. \rangle$ in SVM papers
- $x \in \mathbb{R}^D$ (but not necessarily), but $\phi(x)$ can be in any space – same, higher, lower or even in an infinite dimensional space

# Non-linearly separable data



Not linearly separable

Same data, but now linearly separable

- Want to map our data to a linearly separable space
- Instead of $x$, want $\phi(x)$, in a separable space ($\phi(x)$ is a feature map)

What if $\phi(x)$ is much higher dimensional? We do not want to learn more parameters and mapping could become very expensive

# Radial Basis Function Kernel (RBF)

Arguably the most popular kernel function ( for Support Vector Machine)

$$K(x_i, x_j) = \exp{-\frac{1}{2\sigma^2} \|x_i - x_j\|^2}$$

$\phi(\boldsymbol{x}) =$?

- It is infinite dimensional and fairly involved, no easy way to actually perform the mapping to this space, but we know what an inner product looks like in it

$\sigma = $ ?

- a hyperparameter
- With a really low sigma the model becomes close to a KNN approach (potentially very expensive)

# Some other kernels

Other kernels exist

- Histogram Intersection Kernel
  - good for histogram features
- String kernels
  - specifically for text and sentence features
- Proximity distribution kernel
- (Spatial) pyramid matching kernel

# Kernel CCA

If we remember CCA it used only inner products in definitions when dealing with data, that means we can again use kernels

$$(w_1^*, w_2^*) = \operatorname*{argmax}_{w_1, w_2} \frac{w_1' \Sigma_{12} w_2}{\sqrt{w_1' \Sigma_{11} w_1 w_2' \Sigma_{22} w_2}} = \operatorname*{argmax}_{w_1' \Sigma_{11} w_1 = w_2' \Sigma_{22} w_2 = 1} w_1' \Sigma_{12} w_2$$

We can now map into a high-dimensional non-linear space instead

$$(\alpha_1^*, \alpha_2^*) = \operatorname*{argmax}_{\alpha_1, \alpha_2} \frac{\alpha_1' K_1 K_2 \alpha_2}{\sqrt{(\alpha_1' K_1^2 \alpha_2)(\alpha_1' K_2^2 \alpha_2)}} = \operatorname*{argmax}_{\alpha_1' K_1^2 \alpha_1 = \alpha_2' K_2^2 \alpha_2 = 1} \alpha_1' K_1 K_2 \alpha_2,$$

[Lai et al. 2000]

Language Technologies Institute

Carnegie Mellon University

# Transformer's Attention Function



**Scale dot-product attention:**

$$\boldsymbol{\alpha} = softmax\left(\frac{\boldsymbol{x_q W_q (x_k W_k)}^T}{\sqrt{d_k}}\right)$$

How can you interpret it as a kernel similarity function?

Language Technologies Institute

Carnegie Mellon University

# Transformer's Attention Function



**Scale dot-product attention:**

$$\alpha = softmax\left(\frac{x_q W_q (x_k W_k)^T}{\sqrt{d_k}}\right)$$

**Kernel-formulated attention:**

$$\alpha = \frac{k(x_q, x_k)}{\sum_{\{x_k'\}} k(x_q, x_k')}$$

What is the impact of the kernel function?

Tsai et al., Transformer Dissection: An Unified Understanding for Transformer's Attention via the Lens of Kernel, EMNLP 2019

# Transformer's Attention Function



$h_1$

$\Sigma$

Scale dot-product attention

$\times$

$\alpha_{2,1}$
$\alpha_{1,1}$

$v_1$

$q_1$    $k_1$

$W_q$    $W_k$    $W_v$

$x_1$

## What is the impact of the kernel function?

| Type | Kernel Form | NMT (BLEU↑) | |
|------|-------------|-------------|---|
| | | Asym. ($W_q \neq W_k$) | Sym. ($W_q = W_k$) |
| Linear | $\langle f_a W_q, f_b W_k \rangle$ | not converge | not converge |
| Polynomial | $\left( \langle f_a W_q, f_b W_k \rangle \right)^2$ | 32.72 | 32.43 |
| Exponential | $\exp\left( \frac{\langle f_a W_q, f_b W_k \rangle}{\sqrt{d_k}} \right)$ | 33.98 | 33.78 |
| RBF | $\exp\left( -\frac{\| f_a W_q - f_b W_k \|^2}{\sqrt{d_k}} \right)$ | **34.26** | 34.14 |

Conventional Transformer →

## What is the best way to integrate the position embedding?

Tsai et al., Transformer Dissection: An Unified Understanding for Transformer's Attention via the Lens of Kernel, EMNLP 2019

Language Technologies Institute

Carnegie Mellon University

# Transformer's Attention Function



What is the best way to integrate the position embedding?

| PE Incorporation | Kernel Form | NMT (BLEU↑) |
|---|---|---|
| Direct-Sum | $k_{\exp}\left(f_q + t_q, f_k + t_k\right)$ | 33.98 |
| Look-up Table | $L_{t_q-t_k, f_q} \cdot k_{\exp}\left(f_q, f_k\right)$ | 34.12 |
| Product Kernel | $k_{\exp}\left(f_q, f_k\right) \cdot k_{f_q}\left(t_q, t_k\right)$ | 33.62 |
| Product Kernel | $k_F\left(f_q, f_k\right) \cdot k_T\left(t_q, t_k\right)$ | **34.71** |

Vaswami et al → Direct-Sum

Transformer XL → Product Kernel

Proposed → Product Kernel

with $k_F(f_q, f_k) = \exp\left(\frac{\langle f_q W_F, f_k W_F\rangle}{\sqrt{d_k}}\right)$  Same weight matrices!

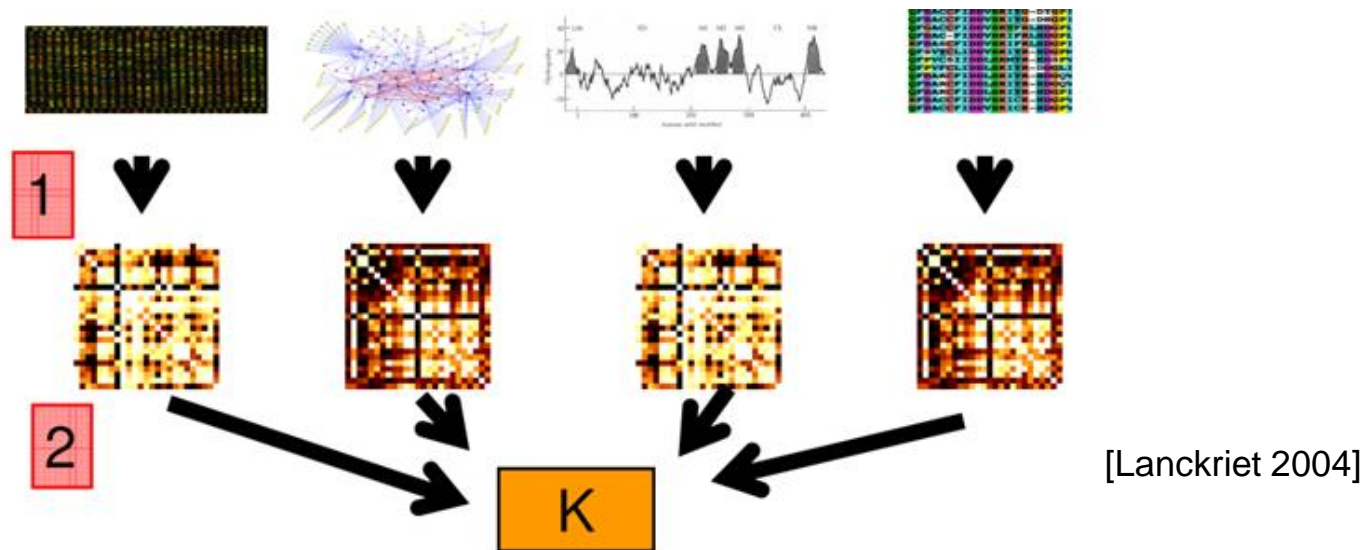and $k_T(t_q, t_k) = \exp\left(\frac{\langle t_q W_T, t_k W_T\rangle}{\sqrt{d_k}}\right),$

Scale dot-product attention

ansformer's

Can Kernels be used as a Fusion Mechanism (for late fusion)?

# Multiple Kernel Learning



[Lanckriet 2004]

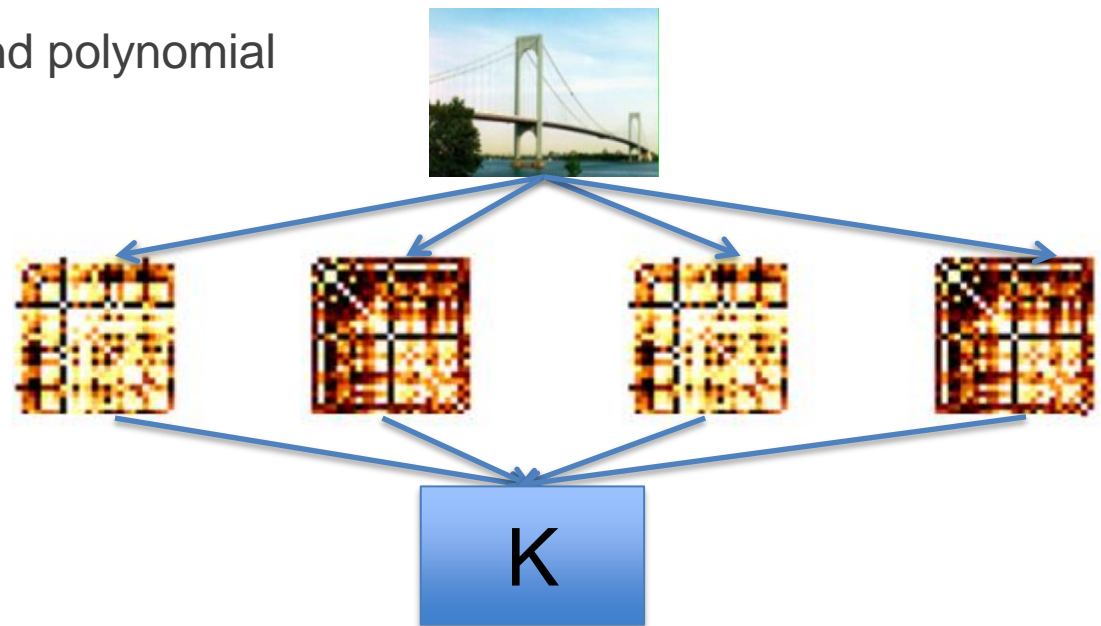- Instead of providing a single kernel and validating which one works optimize in a family of kernels (or different families for different modalities)
- Works well for unimodal and multimodal data, very little adaptation is needed

# MKL in Unimodal Case

- Pick a family of kernels and learn which kernels are important for the classification case

- For example a set of RBF and polynomial kernels

# MKL in Multimodal/Multiview Case

- Pick a family of kernels for each modality and learn which kernels are important for the classification case

- Does not need to be different modalities, often we use different views of the same modality (HOG, SIFT, etc.)

# Co-Learning
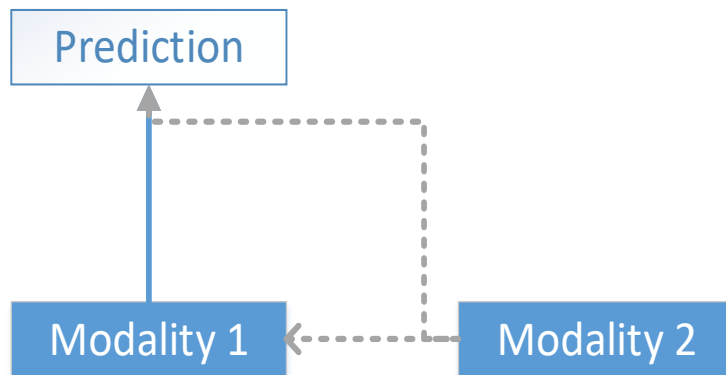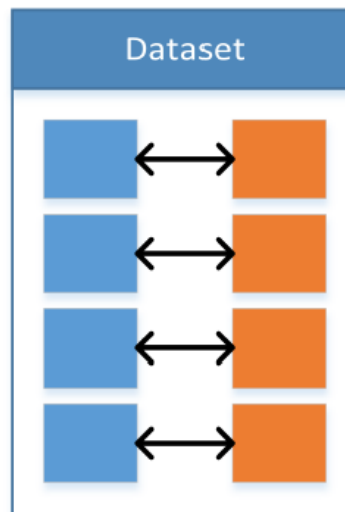
Carnegie Mellon University

# Co-Learning - The 5ᵗʰ Multimodal Challenge

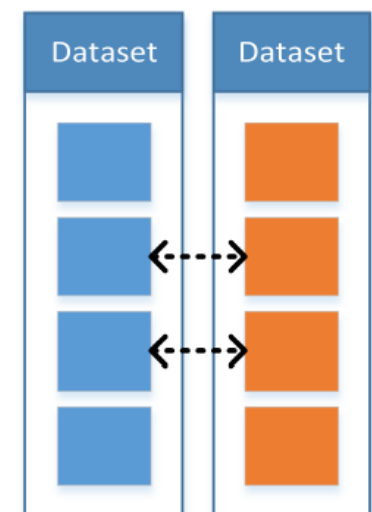**Definition:** Transfer knowledge between modalities, including their representations and predictive models.

Language Technologies Institute

Carnegie Mellon University

# Co-learning Example with Paired Data

Learn vector representations for text using visual co-occurrences

Four types of co-occurrences:

    (a) Object - Attribute
    (b) Attribute - Attribute
    (c) Context
    (d) Object-Hypernym



| Region | Object Words | Attribute Words |
|---|---|---|
| (green) | man, person, adult, mammal | muscular, smiling |
| (blue) | woman, person, adult, mammal | lean, smiling |
| (orange) | table, tablecloth, furniture | striped, oval |
| (dark blue) | rice, carbohydrates, food | white, grainy, cooked |
| (purple) | salad, roughage, food | leafy, chopped, healthy, red, green |
| (red) | glass, glassware, utensil | clear, transparent, reflective, tall |
| (yellow) | plate, crockery, utensil | ceramic, white, round, circular |
| (cyan) | fork, cutlery, utensil | metallic, shiny, reflective |
| (magenta) | spoon, cutlery, utensil | serving, metallic, shiny, reflective |

ViCo: Word Embeddings from Visual Co-occurrences

# ViCo: Word Embeddings from Visual Co-occurrences

## Relatedness through Co-occurrences

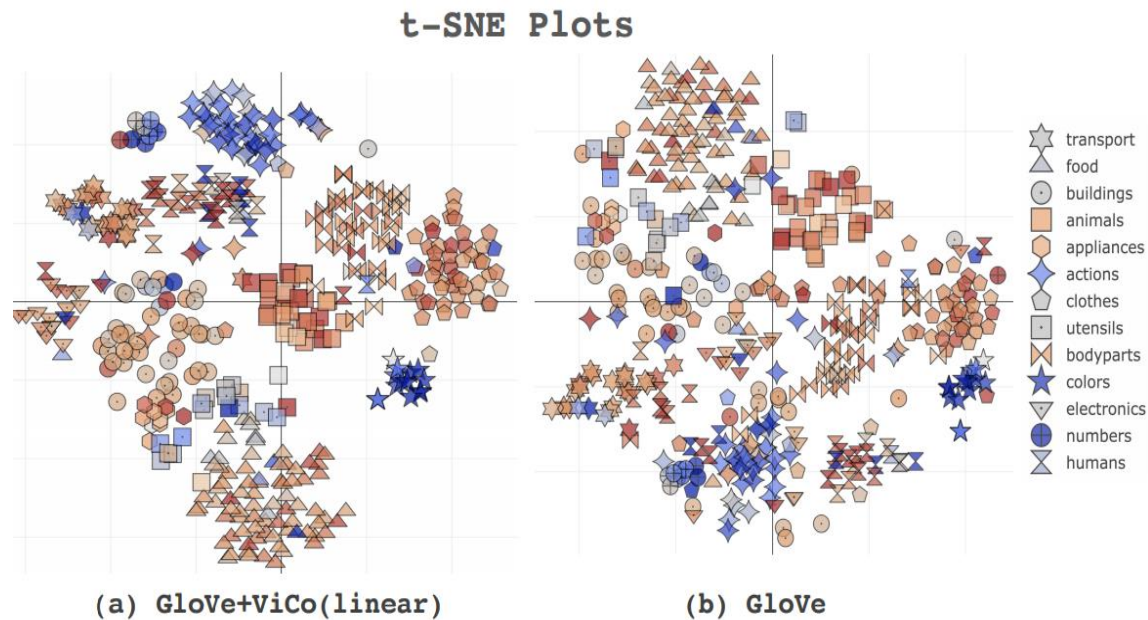| Word Pair | ViCo | Obj-Attr | Attr-Attr | Obj-Hyp | Context | GloVe |
|---|---|---|---|---|---|---|
| crouch / squat | 0.61 | 0.74 | 0.72 | 0.18 | 0.25 | 0.05 |
| sweet / dessert | 0.66 | 0.78 | 0.76 | 0.56 | 0.79 | 0.43 |
| man / male | 0.71 | 0.98 | 0.8 | 0.38 | 1 | 0.34 |
| purple / violet | 0.75 | 0.93 | 1 | 0.24 | 0.03 | 0.52 |
| hosiery / sock | 0.52 | 0.27 | 0.18 | 0.87 | 0.07 | 0.23 |
| aeroplane / aircraft | 0.73 | 0.43 | 0.07 | 0.87 | 0.75 | 0.43 |
| bench / pew | 0.63 | 0.67 | 0.09 | 0.79 | -0.14 | 0.1 |
| keyboard / mouse | 0.19 | 0.63 | 0.19 | 0.09 | 0.95 | 0.52 |
| laptop / desk | 0.39 | 0.23 | 0.24 | 0.1 | 0.94 | 0.28 |
| window / door | 0.59 | 0.46 | 0.35 | 0.53 | 0.93 | 0.67 |
| hair / blonde | 0.16 | 0.56 | 0.32 | -0.15 | 0.17 | 0.51 |
| thigh / ankle | 0.09 | 0.19 | 0.03 | 0.01 | 0.39 | 0.74 |
| garlic / onion | 0.36 | -0.03 | 0.3 | 0.37 | 0.56 | 0.77 |
| driver / car | 0.27 | 0.16 | 0.26 | 0.12 | 0.53 | 0.71 |
| girl / boy | 0.41 | 0.38 | 0.22 | 0.44 | 0.74 | 0.83 |

Since ViCo is learned from multiple types of co-occurrences, it is hypothesized to provide a richer sense of relatedness

➢ Learned using a multi-task Log-Bilinear Model

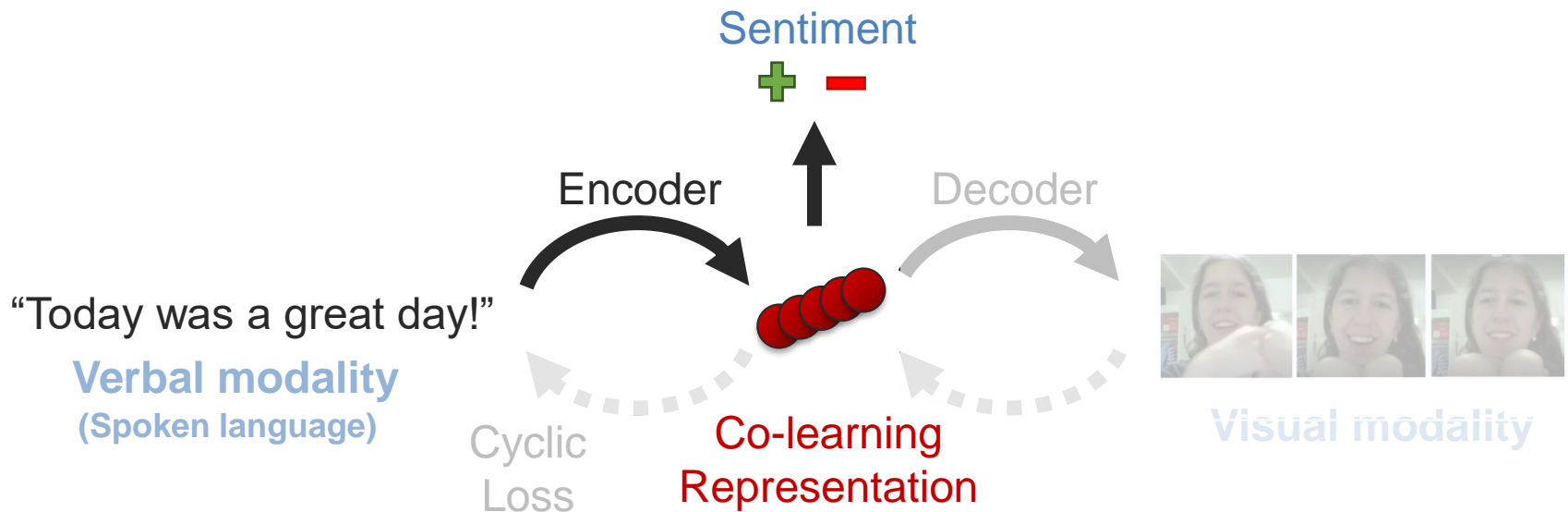# ViCo: Word Embeddings from Visual Co-occurrences
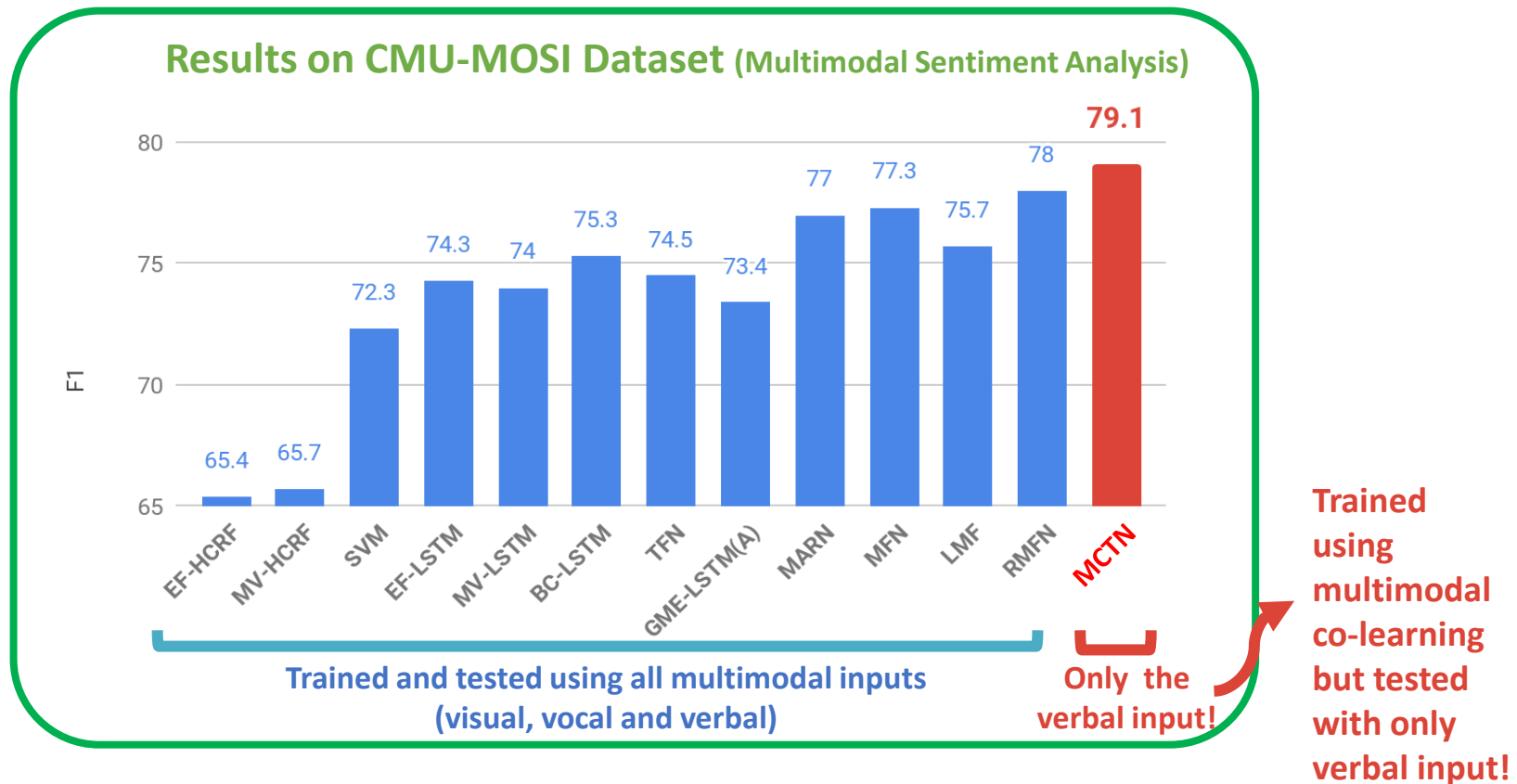
ViCO leads to more homogenous clusters compared to GloVe



t-SNE Plots

(a) GloVe+ViCo(linear)    (b) GloVe

Language Technologies Institute

Carnegie Mellon University

# Another Example of Co-Learning with Paired Data: Multimodal Cyclic Translation



Sentiment

**+** **−**

Encoder

Decoder

"Today was a great day!"

**Verbal modality**
**(Spoken language)**

Cyclic Loss

Co-learning Representation

Visual modality

Paul Pu Liang*, Hai Pham*, et al., "Found in Translation: Learning Robust Joint Representations by Cyclic Translations Between Modalities", AAAI 2019

Carnegie Mellon University

# Another Example of Co-Learning with Paired Data: Multimodal Cyclic Translation



**Results on CMU-MOSI Dataset** (Multimodal Sentiment Analysis)

Trained and tested using all multimodal inputs (visual, vocal and verbal)

Only the verbal input!

Trained using multimodal co-learning but tested with only verbal input!

Paul Pu Liang*, Hai Pham*, et al., "Found in Translation: Learning Robust Joint Representations by Cyclic Translations Between Modalities", AAAI 2019
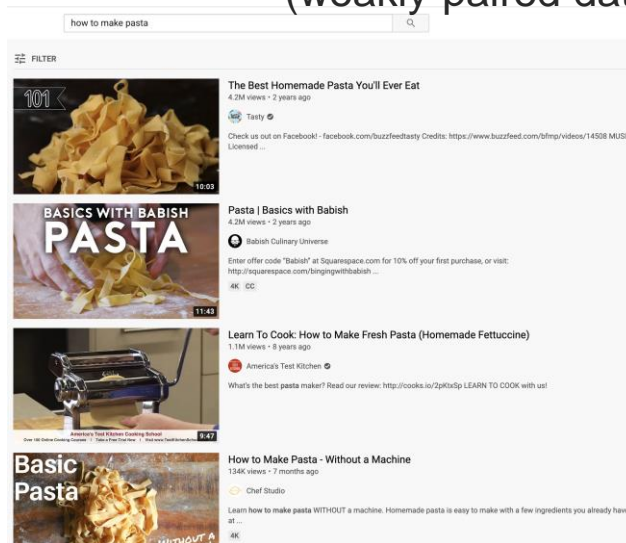
# Co-Learning Example with Weakly Paired Data

*NEW paper*

**Goal:** Learn better visual representations…

… by taking advantage of large-scale video+language resources

### Instructional videos
(weakly-paired data)



*it's turning into a much thicker mixture*
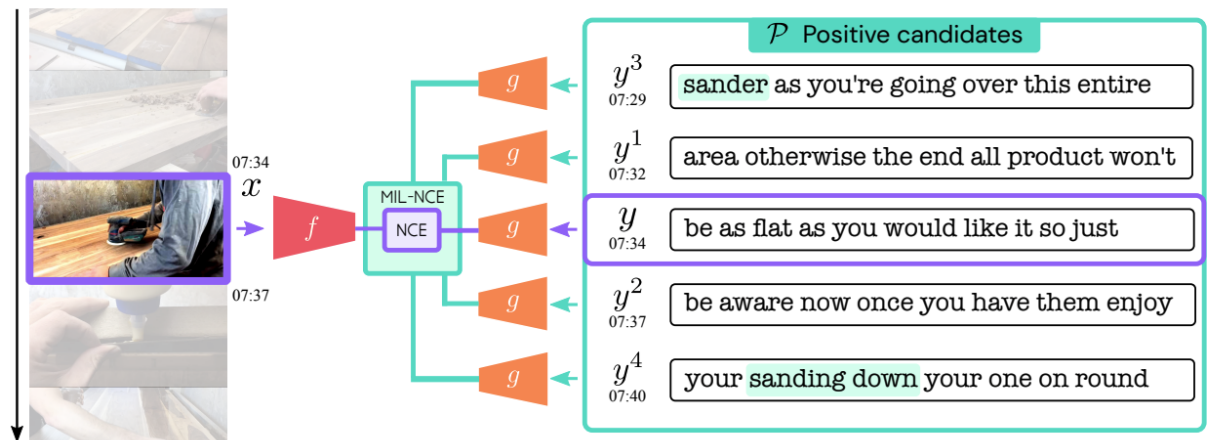


*The biggest mistake is not kneading it enough*



…

End-to-End Learning of Visual Representations from Uncurated Instructional Videos
Antoine Miech, Jean-Baptiste Alayrac, Lucas Smaira, Ivan Laptev, Josef Sivic, and Andrew Zisserman – CVPR 2020

# Weakly Paired Data

**Data point:** "a short 3.2 seconds video clip (32 frames at 10 FPS) together with a small number of words (not exceeding 16)"



## How to handle this misalignment? Multi-instance learning!

## How to do it self-supervised? Contrastive learning!

End-to-End Learning of Visual Representations from Uncurated Instructional Videos
Antoine Miech, Jean-Baptiste Alayrac, Lucas Smaira, Ivan Laptev, Josef Sivic, and Andrew Zisserman – CVPR 2020

Language Technologies Institute

Carnegie Mellon University

# Multiple Instance Learning Noise Contrastive Estimation

**Objective**

Given video $x$ and text $y$ from a positive set $P_i$ and a negative set $N_i$, maximize the positive / total score ratio

$$\max_{f,g} \sum_{i=1}^{n} \log \left( \frac{\sum\limits_{(x,y)\in\mathcal{P}_i} e^{f(x)^\top g(y)}}{\sum\limits_{(x,y)\in\mathcal{P}_i} e^{f(x)^\top g(y)} + \sum\limits_{(x',y')\sim\mathcal{N}_i} e^{f(x')^\top g(y')}} \right)$$

Note: Doing so requires maximizing $f(x)^\top g(y)$ for only positive examples

1. Using sets of positive and negative examples to ~wash out the misaligned text
2. Ideally, we would maximize all positives over all possible negatives (intractable)

Language Technologies Institute

**Carnegie Mellon University**

# Experiments – HowTo100M Dataset



End-to-End Learning of Visual Representations from Uncurated Instructional Videos
Antoine Miech, Jean-Baptiste Alayrac, Lucas Smaira, Ivan Laptev, Josef Sivic, and Andrew Zisserman – CVPR 2020

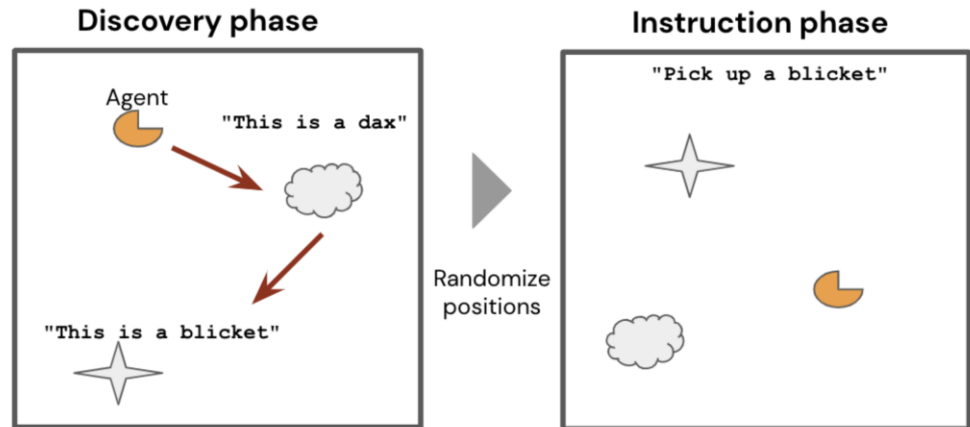# Research Trend: Few-Shot Learning and Weakly Supervised

# Few-Shot Learning in RL Environment

**Discovery** phase:
- Explore environment and when the agent sees an object, a description is provided to it.

**Instruction** phase:
- Given an instruction, e.g. "Pick up a dax",+1.0 reward if picked up correct object

**One-shot:** never seen "theble"

➡ "Fast-mapping"

**Key idea:** Dual-coding Episodic Memory architecture
*(a slow one and a fast one)*



Hill et al., Grounded Language Learning Fast and Slow. arXiv 2020
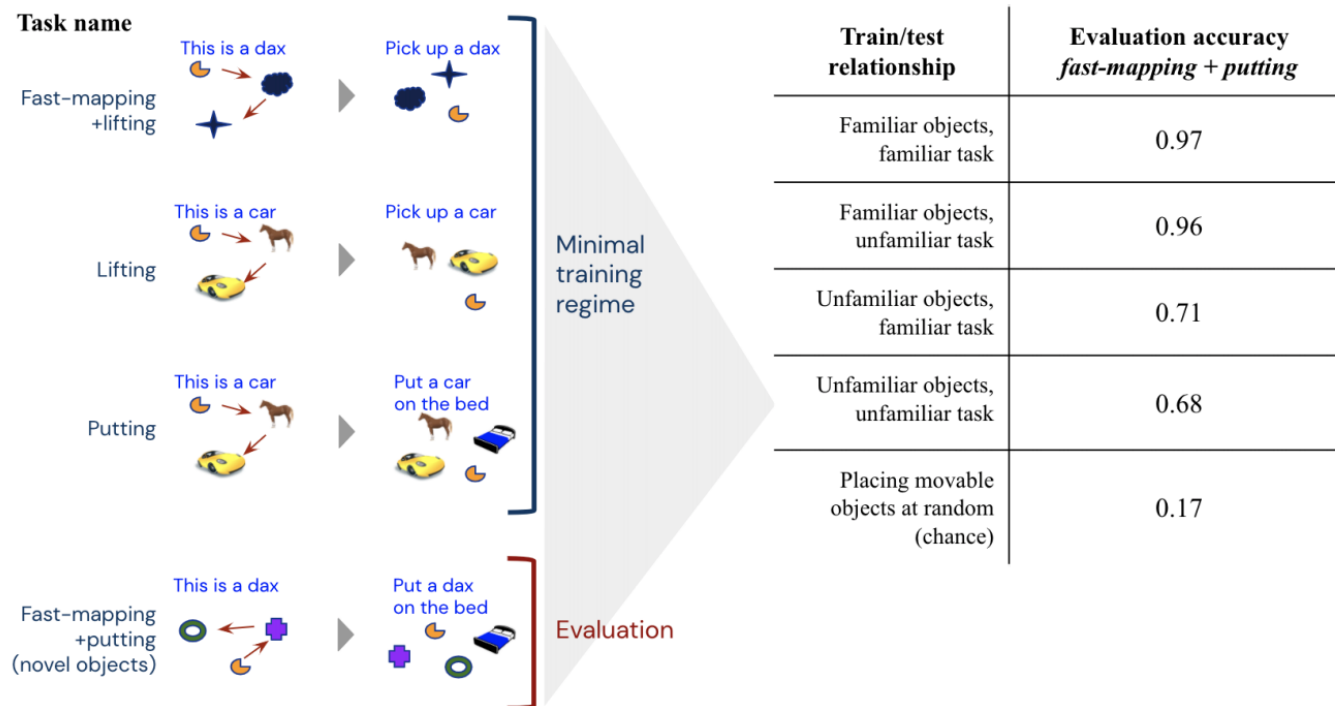
# Grounded Language Learning

Generalization to new objects and new instructions.



Hill et al., Grounded Language Learning Fast and Slow. arXiv 2020

# Grounded Language Learning

Generalization to new objects and new instructions.



Hill et al., Grounded Language Learning Fast and Slow. arXiv 2020

# Weakly-Supervised Phrase Grounding

**NEW** paper

**Phrase grounding** is a task that studies the mapping from textual phrases to regions of an image. But limited labeled data…



**Supervised**
phrase-object annotations

**Weakly-supervised**
image-caption annotations

A young baby crawls across the wood floor towards the water bottle

A young baby crawls across the wood floor towards the water bottle

**General solution:** leverage more caption-image datasets, which can then be used as a form of weak supervision

MAF: Multimodal Alignment Framework for Weakly-Supervised Phrase Grounding, EMNLP 2020

Language Technologies Institute

Carnegie Mellon University

# Multimodal Alignment Framework



An older gentleman is standing next to the man with a red accordion over his shoulder.

localization with find-grained visual representation

**Specific solution:**
Enhance visual representations of objects (e.g., man) by "shifting" it based on the caption phrases.
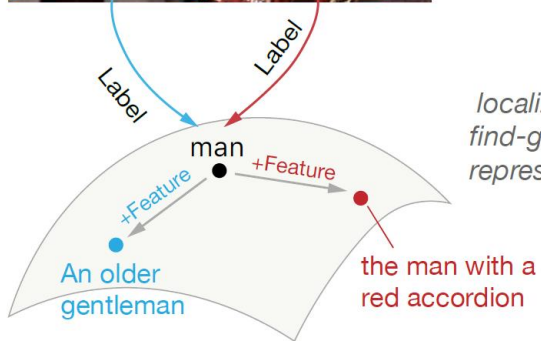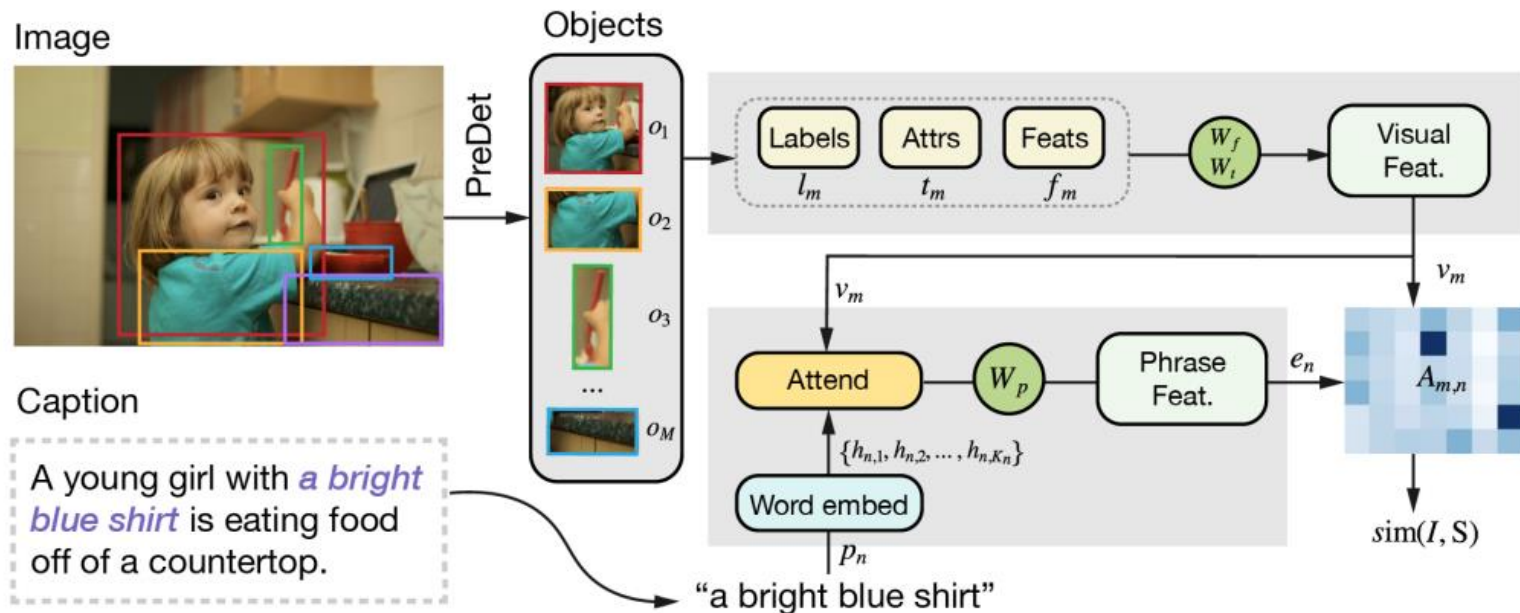
➡️ Fine-grained visual representations

How?

Contrastive learning!

MAF: Multimodal Alignment Framework for Weakly-Supervised Phrase Grounding, EMNLP 2020

# Multimodal Alignment Framework



Image

Caption

A young girl with *a bright blue shirt* is eating food off of a countertop.

Objects

PreDet

Labels $l_m$  Attrs $t_m$  Feats $f_m$  $W_f$ $W_t$  Visual Feat.

$v_m$

Attend  $W_p$  Phrase Feat.  $e_n$  $A_{m,n}$

$\{h_{n,1}, h_{n,2}, ..., h_{n,K_n}\}$

Word embed

$p_n$

"a bright blue shirt"

$sim(I, S)$

Contrastive loss: $\mathcal{L} = -\log \dfrac{e^{\text{sim}(I,S)}}{\sum_{I' \in batch} e^{\text{sim}(I',S)}}.$

MAF: Multimodal Alignment Framework for Weakly-Supervised Phrase Grounding, EMNLP 2020

Language Technologies Institute

Carnegie Mellon University

# MAF: Multimodal Alignment Framework for Weakly-Supervised Phrase Grounding

| Method | Vis. Features | Acc. (%) | UB |
|---|---|---|---|
| **Supervised** | | | |
| GroundeR (Rohrbach et al., 2016) | $VGG_{det}$ | 47.81 | 77.90 |
| CCA (Plummer et al., 2015) | $VGG_{det}$ | 50.89 | 85.12 |
| BAN (Kim et al., 2018) | ResNet-101 | 69.69 | 87.45 |
| visualBERT (Li et al., 2019) | ResNet-101 | 71.33 | 87.45 |
| DDPN (Yu et al., 2018) | ResNet-101 | 73.30 | - |
| CGN (Liu et al., 2020) | ResNet-101 | 76.74 | - |
| **Weakly-Supervised** | | | |
| GroundeR (Rohrbach et al., 2016) | $VGG_{det}$ | 28.93 | 77.90 |
| Link (Yeh et al., 2018) | $YOLO_{det}$ | 36.93 | - |
| KAC (Chen et al., 2018) | $VGG_{det}$ | 38.71 | - |
| MAF (Ours) | $VGG_{det}$ | 44.39 | 86.29 |
| MAF (Ours) | ResNet-101 | **61.43** | 86.29 |

MAF: Multimodal Alignment Framework for Weakly-Supervised Phrase Grounding, EMNLP 2020

Language Technologies Institute

Carnegie Mellon University

# References

Carnegie Mellon University

# Few-Shot and Weakly Supervised Learning

- [MFAS: Multimodal Fusion Architecture Search](), CVPR 2019

# Few-Shot and Weakly Supervised Learning

- [Grounded Language Learning Fast and Slow](). arxiv 2020
- [MAF: Multimodal Alignment Framework for Weakly-Supervised Phrase Grounding]() EMNLP 2020