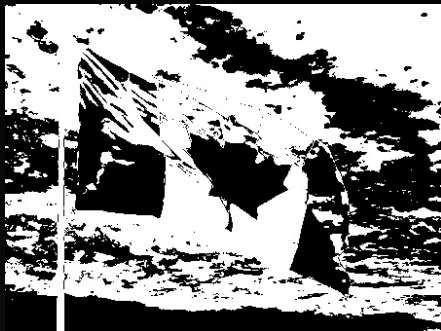


# From language to vision and back again



Center for Brains,  
Minds & Machines

Andrei Barbu

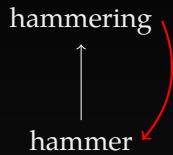


© Source Unknown. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>.









Perception is unreliable.

Top-down knowledge affects our perception.

One integrated representation for many tasks.

Recognition

Retrieval

Generation

Question answering

Disambiguation

Acquiring language

Images, not videos

Translation

Planning

Theory of mind

...



# Humans perform language–vision tasks all the time

Give me the cup.

Which chair should I sit in?

This is an apple.

To win this game you have to make a straight line out of your pieces.

Recognition

S(sentence, video)

Retrieval

Generation

Question answering

Disambiguation

Acquiring language

Images, not videos

Translation

Planning

Theory of mind

...



The **person** **rode** the **skateboard** leftward.

object detector, tracker, **event recognizer**



# Object detection

Figure removed due to copyright restrictions. Please see the video.

Source: Barbu, Andrei, Aaron Michaux, Siddharth Narayanaswamy, and Jeffrey Mark Siskind. "Simultaneous object detection, tracking, and event recognition." *Advances in Cognitive Systems*: 203-220 (2012).

# Object detection

Figure removed due to copyright restrictions. Please see the video.

Source: Barbu, Andrei, Aaron Michaux, Siddharth Narayanaswamy, and Jeffrey Mark Siskind. "Simultaneous object detection, tracking, and event recognition." *Advances in Cognitive Systems*: 203-220 (2012).

# Object detection

Figure removed due to copyright restrictions. Please see the video.

Source: Barbu, Andrei, Aaron Michaux, Siddharth Narayanaswamy, and Jeffrey Mark Siskind. "Simultaneous object detection, tracking, and event recognition." *Advances in Cognitive Systems*: 203-220 (2012).

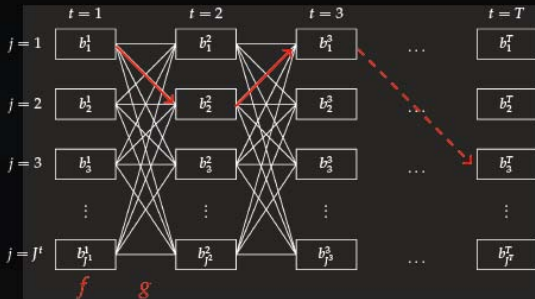
# Object detectors work poorly

Figure removed due to copyright restrictions. Please see the video.

Source: Russakovsky, Olga, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang et al. "Imagenet large scale visual recognition challenge." *International Journal of Computer Vision* 115, no. 3 (2015): 211-252.



# Fixing bad detectors with higher-level knowledge



detection / object / frame  
temporally coherent track  
object detector confidence ( $f$ )  
motion coherence ( $g$ )  
optimal path through the  
lattice of detections  
dynamic programming  
Bellman (1957), Viterbi (1967)

$$\max_{j^1, \dots, j^T} \sum_{t=1}^T f(b_{j^t}^t) + \sum_{t=2}^T g(b_{j^{t-1}}^{t-1}, b_{j^t}^t)$$

# Fixing bad detectors with higher-level knowledge

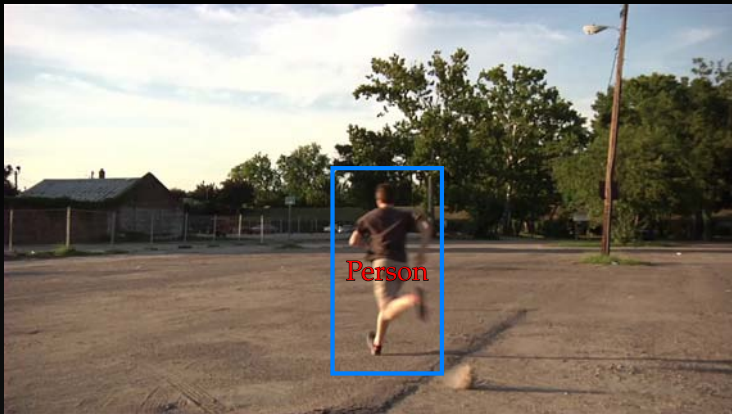


Courtesy of Andrei Barbu, Alexander Bridge, Dan Coroian, Sven Dickinson, Sam Mussman, Siddharth Narayanaswamy, Dhaval Salvi, Lara Schmidt, Jiangnan Shangguan, Jeffrey Mark Siskind, Jarrell Waggoner, Song Want, Jinlian Wei, Yifan Yin & Zhiqi Zhang. Used with permission.  
Source: Barbu, Andrei, Alexander Bridge, Dan Coroian, Sven Dickinson, Sam Mussman, Siddharth Narayanaswamy, Dhaval Salvi et al. "Large-scale automatic labeling of video events with verbs based on event-participant interaction." arXiv preprint arXiv:1204.3616 (2012).

detection / object / frame  
temporally coherent track  
object detector confidence ( $f$ )  
motion coherence ( $g$ )  
optimal path through the  
lattice of detections  
dynamic programming  
Bellman (1957), Viterbi (1967)

$$\max_{j^1, \dots, j^T} \sum_{t=1}^T f(b_{j^t}^t) + \sum_{t=2}^T g(b_{j^{t-1}}^{t-1}, b_{j^t}^t)$$

# Feature vector — single participant



position	$\frac{d}{dt}$ position	$\frac{d^2}{dt^2}$ position	
aspect ratio	$\frac{d}{dt}$ aspect ratio	area	$\frac{d}{dt}$ area
<b>object class</b>	root-filter index		

# Feature vector — dual participant



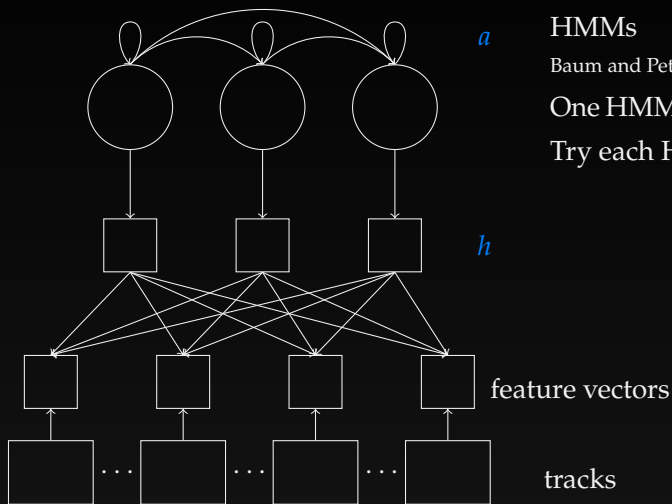
distance     $\frac{d}{dt}$ distance    orientation     $\frac{d}{dt}$ orientation

person riding skateboard

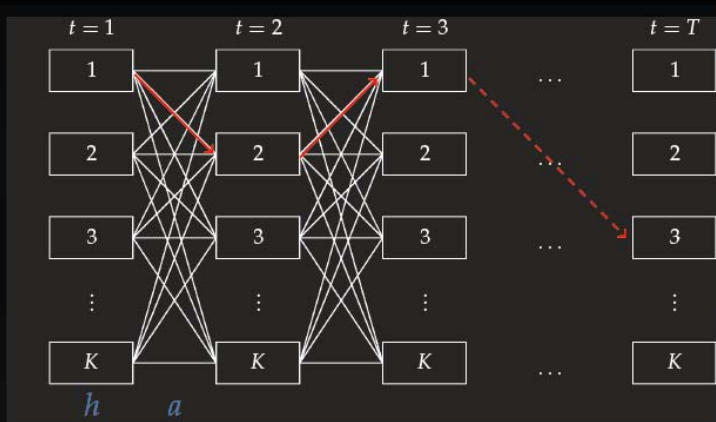
person approaching person

skateboard approaching person

# Event recognition



# Event recognition



$$\max_{k^1, \dots, k^T} \sum_{t=1}^T h(k^t, b_{j^t}^t) + \sum_{t=2}^T a(k^{t-1}, k^t)$$





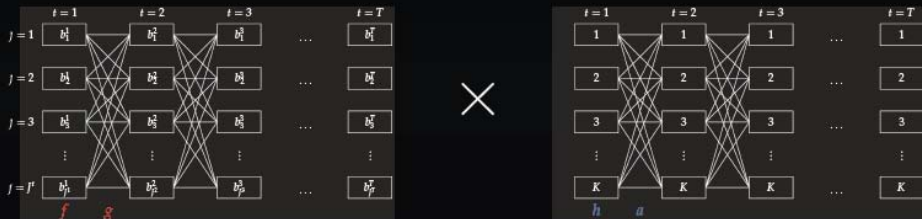


# Tracking in the context of event recognition



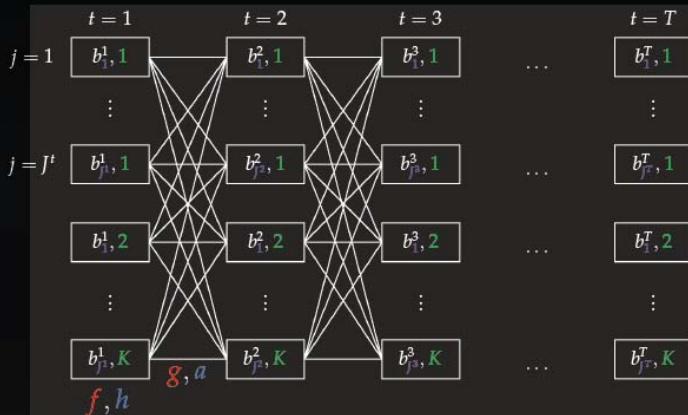
$$\max_{j^1, \dots, j^T} \sum_{t=1}^T f(b_{j^t}^t) + \sum_{t=2}^T g(b_{j^{t-1}}^{t-1}, b_{j^t}^t) + \max_{k^1, \dots, k^T} \sum_{t=1}^T h(k^t, b_{j^t}^t) + \sum_{t=2}^T a(k^{t-1}, k^t)$$

# Tracking in the context of event recognition



$$\max_{j^1, \dots, j^T} \max_{k^1, \dots, k^T} \sum_{t=1}^T f(b_{j^t}^t) + \sum_{t=2}^T g(b_{j^{t-1}}^{t-1}, b_{j^t}^t) + \sum_{t=1}^T h(k^t, b_{j^t}^t) + \sum_{t=2}^T a(k^{t-1}, k^t)$$

# Tracking in the context of event recognition



$$\max_{j^1, \dots, j^T} \max_{k^1, \dots, k^T} \sum_{t=1}^T f(b_{j^t}^t) + \sum_{t=2}^T g(b_{j^{t-1}}^{t-1}, b_{j^t}^t) + \sum_{t=1}^T h(k^t, b_{j^t}^t) + \sum_{t=2}^T a(k^{t-1}, k^t)$$

# Tracking in the context of event recognition in action



tracking

tracking and event recognition

Courtesy of Andrei Barbu, Alexander Bridge, Dan Coroian, Sven Dickinson, Sam Mussman, Siddharth Narayanaswamy, Dhaval Salvi, Lara Schmidt, Jiangnan Shangguan, Jeffrey Mark Siskind, Jarrell Waggoner, Song Want, Jinlian Wei, Yifan Yin & Zhiqi Zhang. Used with permission.  
Source: Barbu, Andrei, Alexander Bridge, Dan Coroian, Sven Dickinson, Sam Mussman, Siddharth Narayanaswamy, Dhaval Salvi et al. "Large-scale automatic labeling of video events with verbs based on event-participant interaction." arXiv preprint arXiv:1204.3616 (2012).

# Tracking in the context of event recognition in action



tracking

tracking and event recognition



# Building sentences out of **trackers** and **words**

## Viterbi tracker

track 1



$$\max_{j_1^1, \dots, j_1^T}$$

$$\sum_{t=1}^T f(b_{j_1^t}^t) + \sum_{t=2}^T g(b_{j_1^{t-1}}^{t-1}, b_{j_1^t}^t)$$

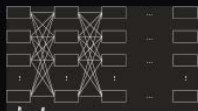
# Building sentences out of **trackers** and **words**

Event tracker for intransitive verbs

**track 1**



×



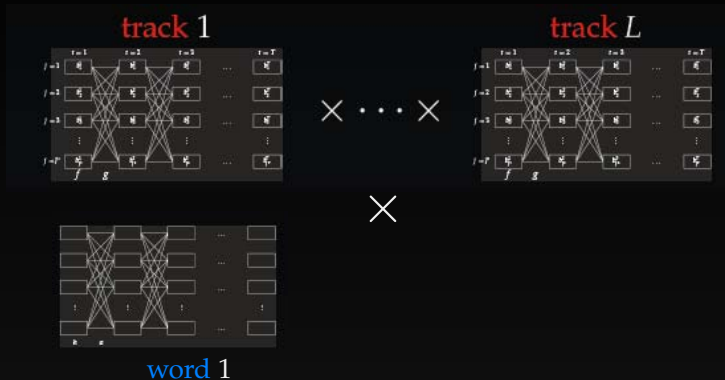
**word 1**

$$\max_{j_1^1, \dots, j_1^T} \max_{k_1^1, \dots, k_1^T} \sum_{t=1}^T f(b_{j_1^t}^t) + \sum_{t=2}^T g(b_{j_1^{t-1}}^{t-1}, b_{j_1^t}^t) + \sum_{t=1}^T h(k^t, b_{j_1^t}^t) + \sum_{t=2}^T a(k^{t-1}, k^t)$$



# Building sentences out of **trackers** and **words**

Event tracker

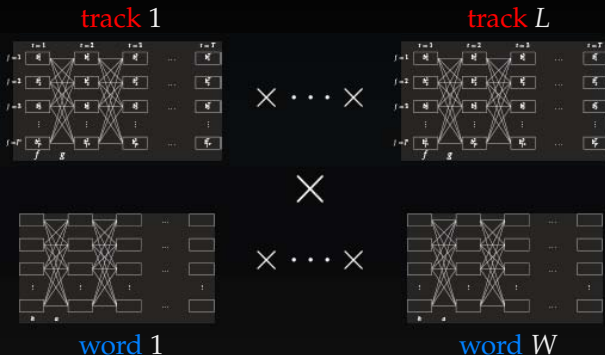


$$\max_{j_1^1, \dots, j_1^T} \max_{k_1^1, \dots, k_1^T} \sum_{l=1}^L \sum_{t=1}^T f(b_{j_l^t}^t) + \sum_{t=2}^T g(b_{j_l^{t-1}}^{t-1}, b_{j_l^t}^t) + \sum_{t=1}^T h(k^t, b_{j_{\theta^1}^t}^t, b_{j_{\theta^2}^t}^t) + \sum_{t=2}^T a(k^{t-1}, k^t)$$

$\vdots$   
 $j_L^1, \dots, j_L^T$

# Building sentences out of **trackers** and **words**

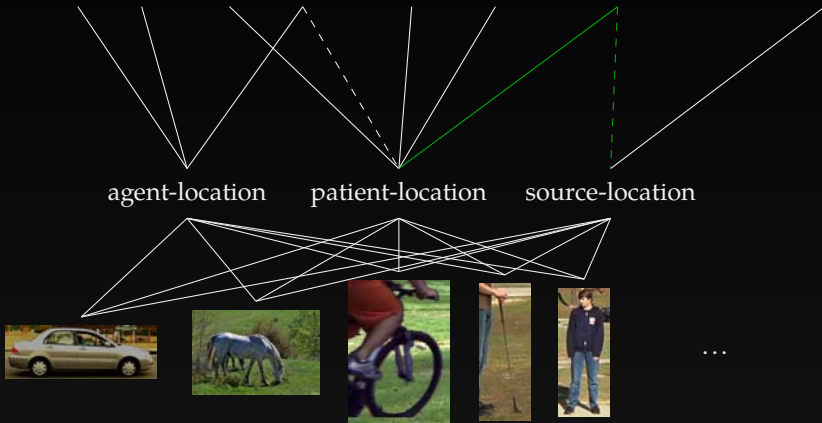
## Sentence tracker



$$\begin{aligned}
 & \max_{j_1^1, \dots, j_1^T} \max_{k_1^1, \dots, k_1^T} \sum_{l=1}^L \sum_{t=1}^T f(b_{j_l^t}^t) + \sum_{t=2}^T g(b_{j_l^{t-1}}^{t-1}, b_{j_l^t}^t) + \sum_{w=1}^W \sum_{t=1}^T h_w(k_w^t, b_{j_{\theta_w^1}^t}^t, b_{j_{\theta_w^2}^t}^t) + \sum_{t=2}^T a_w(k_w^{t-1}, k_w^t) \\
 & \vdots \\
 & \max_{j_L^1, \dots, j_L^T} \max_{k_W^1, \dots, k_W^T} \dots
 \end{aligned}$$

# Sentences

The tall person quickly rode the horse leftward **away from** the other horse.



© Source Unknown. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>.

# Understanding sentences as a whole



© Journal of Artificial Intelligence Research. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>.  
Source: Yu, Haonan, N. Siddharth, Andrei Barbu, and Jeffrey Mark Siskind. "A Compositional Framework for Grounding Language Inference, Generation, and Acquisition in Video." J. Artif. Intell. Res.(JAIR) 52 (2015): 601-713.

We can differentiate events based on the **verb**:

*The person **picked up** an object.*

*The person **put down** an object.*

# Understanding sentences as a whole



© Journal of Artificial Intelligence Research. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>.  
Source: Yu, Haonan, N. Siddharth, Andrei Barbu, and Jeffrey Mark Siskind. "A Compositional Framework for Grounding Language Inference, Generation, and Acquisition in Video." J. Artif. Intell. Res.(JAIR) 52 (2015): 601-713.

We can differentiate events based on the **subject noun**:

The *backpack* approached the bin.

The *chair* approached the bin.

# Understanding sentences as a whole



© Journal of Artificial Intelligence Research. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>.  
Source: Yu, Haonan, N. Siddharth, Andrei Barbu, and Jeffrey Mark Siskind. "A Compositional Framework for Grounding Language Inference, Generation, and Acquisition in Video."  
J. Artif. Intell. Res.(JAIR) 52 (2015): 601-713.

We can differentiate events based on an **adjective in the subject NP**:

*The red object approached the chair.*

*The blue object approached the chair.*

# Understanding sentences as a whole



© Journal of Artificial Intelligence Research. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>.

Source: Yu, Haonan, N. Siddharth, Andrei Barbu, and Jeffrey Mark Siskind. "A Compositional Framework for Grounding Language Inference, Generation, and Acquisition in Video."

J. Artif. Intell. Res.(JAIR) 52 (2015): 601-713.

We can differentiate events based on a **preposition in the object NP**:

*The person picked up an object **to the right of** the bin.*

*The person picked up an object **to the left of** the bin.*

Siddharth et al 2014

August 28, 2015

Recognition

$S(\text{sentence}, \text{video})$

Retrieval

$\operatorname{argmax}_{v \in V} S(s, v)$

Generation

Question answering

Disambiguation

Acquiring language

Images, not videos

Translation

Planning

Theory of mind

...



# Sentential retrieval



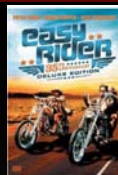
© Warner Bros. Family Entertainment. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>.



© United Artists. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>.



© Warner Bros. Family Entertainment. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>.



© Columbia Pictures. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>.



© Buena Vista Pictures. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>.



© Metro-Goldwyn-Mayer. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>.



© Columbia Pictures. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>.



© Universal Pictures. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>.



© United Artists. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>.



© Warner Bros. Family Entertainment. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>.

Recognition

$S(\text{sentence}, \text{video})$

Retrieval

$\operatorname{argmax}_{v \in V} S(s, v)$

Generation

$\operatorname{argmax}_{s \in L} S(s, v)$

Question answering

Disambiguation

Acquiring language

Images, not videos

Translation

Planning

Theory of mind

...

# Generating sentences

$S \rightarrow NP VP$   
 $NP \rightarrow D [A] N [PP]$   
 $D \rightarrow an \mid the$   
 $A \rightarrow blue \mid red$   
 $N \rightarrow person \mid backpack \mid chair \mid bin \mid object$   
 $PP \rightarrow P NP$   
 $P \rightarrow to\ the\ left\ of \mid to\ the\ right\ of$   
 $VP \rightarrow V NP [Adv] [PP_M]$   
 $V \rightarrow approached \mid carried \mid picked\ up \mid put\ down$   
 $Adv \rightarrow quickly \mid slowly$   
 $PP_M \rightarrow P_M NP$   
 $P_M \rightarrow towards \mid away\ from$

147,123,874,800 sentences without recursion

“the person carried the backpack”

# Generated sentences



The person to the right of the bin picked up the backpack.

© Journal of Artificial Intelligence Research. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>.

Source: Yu, Haonan, N. Siddharth, Andrei Barbu, and Jeffrey Mark Siskind. "A Compositional Framework for Grounding Language Inference, Generation, and Acquisition in Video."

J. Artif. Intell. Res.(JAIR) 52 (2015): 601-713.

Recognition

$$S(\text{sentence}, \text{video})$$

Retrieval

$$\operatorname{argmax}_{v \in V} S(s, v)$$

Generation

$$\operatorname{argmax}_{s \in L} S(s, v)$$

Question answering

$$\operatorname{argmax}_{s \in L} S'(s, s_q, v)$$

Disambiguation

Acquiring language

Images, not videos

Translation

Planning

Theory of mind

...

Recognition

$$S(\text{sentence}, \text{video})$$

Retrieval

$$\operatorname{argmax}_{v \in V} S(s, v)$$

Generation

$$\operatorname{argmax}_{s \in L} S(s, v)$$

Question answering

$$\operatorname{argmax}_{s \in L} S(Q(s, s_q), v)$$

Disambiguation

Acquiring language

Images, not videos

Translation

Planning

Theory of mind

...

# Question answering



What did the person put on top of the red car?

The person put **NP** on top of the red car.

The person put **the pear** on top of the red car.

# Question answering



Who put an object on top of the red car?



# Generation for question answering

Who put an object on top of the red car?

NP put an object on top of the red car.

...

The person put an object on top of the red car.

# Generation for question answering

Who put an object on top of the red car?

NP put an object on top of the red car.

...

~~The person~~ put an object on top of the red car.

# *Discriminative* generation for question answering

Who put an object on top of the red car?

NP put an object on top of the red car.

...

~~The person~~ put an object on top of the red car.

...

The person on the left of the car put an object on top of the red car.

# Question answering



Who put an object on top of the red car?

NP put an object on top of the red car.

The person on the left of the car put an object on top of the red car.

Recognition

$$S(\text{sentence}, \text{video})$$

Retrieval

$$\operatorname{argmax}_{v \in V} S(s, v)$$

Generation

$$\operatorname{argmax}_{s \in L} S(s, v)$$

Question answering

$$\operatorname{argmax}_{s \in L} S(Q(s, s_q), v)$$

Disambiguation

$$\operatorname{argmax}_{i \in \text{parser}(s)} S(i, v)$$

Acquiring language

Images, not videos

Translation

Planning

Theory of mind

...

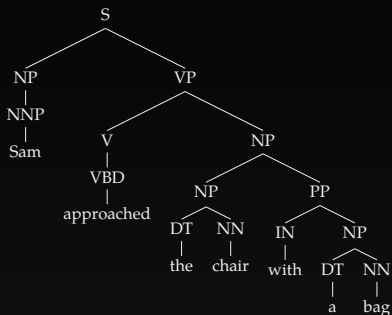
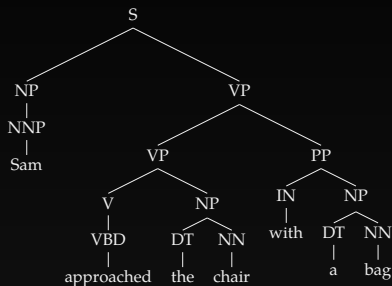
# Disambiguation

Danny approached the chair with a bag.



# Disambiguation

Danny approached the chair with a bag.



## PP Attachment

Danny looked at Andrei with a telescope.



Courtesy of Yevgeni Berzak, Andrei Barbu, Daniel Harari, Boris Katz & Shimon Ullman. License CC BY.  
Source: Berzak, Yevgeni, Andrei Barbu, Daniel Harari, Boris Katz, and Shimon Ullman. "Do you see what i mean? visual resolution of linguistic ambiguities." arXiv preprint arXiv:1603.08079 (2016).



# Disambiguation

PP Attachment

VP Attachment

Andrei approached the person holding a green chair.



# Disambiguation

PP Attachment

VP Attachment

Conjunction

Danny and Andrei picked up the yellow bag and chair.



Courtesy of Yevgeni Berzak, Andrei Barbu, Daniel Harari, Boris Katz & Shimon Ullman. License CC BY.  
Source: Berzak, Yevgeni, Andrei Barbu, Daniel Harari, Boris Katz, and Shimon Ullman. "Do you see what i mean? visual resolution of linguistic ambiguities." arXiv preprint arXiv:1603.08079 (2016).

# Disambiguation

PP Attachment

VP Attachment

Conjunction

Logical Form

Someone put down the bags.



# Disambiguation

PP Attachment

VP Attachment

Conjunction

Logical Form

Anaphora

Danny picked up the bag and the chair. It is yellow.



Courtesy of Yevgeni Berzak, Andrei Barbu, Daniel Harari, Boris Katz & Shimon Ullman. License CC BY.  
Source: Berzak, Yevgeni, Andrei Barbu, Daniel Harari, Boris Katz, and Shimon Ullman. "Do you see what i mean? visual resolution of linguistic ambiguities." arXiv preprint arXiv:1603.08079 (2016).

# Disambiguation

PP Attachment

VP Attachment

Conjunction

Logical Form

Anaphora

Ellipsis

Danny left Andrei. Also Yevgeni.



Courtesy of Yevgeni Berzak, Andrei Barbu, Daniel Harari, Boris Katz & Shimon Ullman. License CC BY.  
Source: Berzak, Yevgeni, Andrei Barbu, Daniel Harari, Boris Katz, and Shimon Ullman. "Do you see what i mean? visual resolution of linguistic ambiguities." arXiv preprint arXiv:1603.08079 (2016).

# Not just parse trees

Danny and Andrei moved a chair.

Danny and Andrei move the **same** chair.

**chair**( $x$ )

**move**(Danny,  $x$ ), **move**(Andrei,  $x$ )

Danny and Andrei move **different** chairs.

**chair**( $x$ ), **chair**( $y$ )

**move**(Danny,  $x$ ), **move**(Andrei,  $y$ ),  $x \neq y$

# Not just parse trees

Danny and Andrei moved a chair.

Danny and Andrei move the **same** chair.

**chair**( $x$ ), **person**( $y$ ), **person**( $z$ ),  $y \neq z$   
**move**( $y, x$ ), **move**( $z, x$ )

Danny and Andrei move **different** chairs.

**chair**( $x$ ), **chair**( $y$ ), **person**( $z$ ), **person**( $w$ ),  $z \neq w$   
**move**( $z, x$ ), **move**( $w, y$ ),  $x \neq y$

Recognition

$$S(\text{sentence}, \text{video})$$

Retrieval

$$\operatorname{argmax}_{v \in V} S(s, v)$$

Generation

$$\operatorname{argmax}_{s \in L} S(s, v)$$

Question answering

$$\operatorname{argmax}_{s \in L} S(Q(s, s_q), v)$$

Disambiguation

$$\operatorname{argmax}_{i \in \text{parser}(s)} S(i, v)$$

Acquiring language

$$\operatorname{argmax}_p \prod_{s, v} S(s(p), v)$$

Images, not videos

Translation

Planning

Theory of mind

...



Split into two tasks:

Learning word meaning

Learning syntax

# Language learning: word meaning



The person picked up the chair.



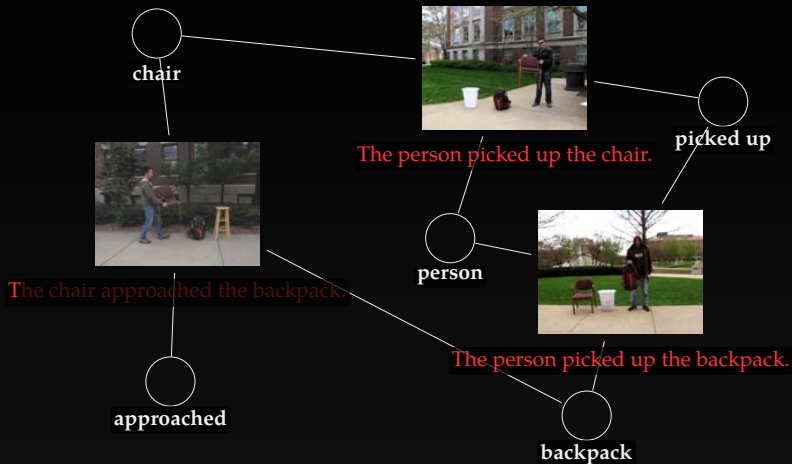
The chair approached the backpack.



The person picked up the backpack.

© Journal of Artificial Intelligence Research. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>.  
Source: Yu, Haonan, N. Siddharth, Andrei Barbu, and Jeffrey Mark Siskind. "A Compositional Framework for Grounding Language Inference, Generation, and Acquisition in Video." J. Artif. Intell. Res.(JAIR) 52 (2015): 601-713.

# Language learning: word meaning



© Journal of Artificial Intelligence Research. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>.  
Source: Yu, Haonan, N. Siddharth, Andrei Barbu, and Jeffrey Mark Siskind. "A Compositional Framework for Grounding Language Inference, Generation, and Acquisition in Video." J. Artif. Intell. Res.(JAIR) 52 (2015): 601-713.

Split into two tasks:

Learning word meaning

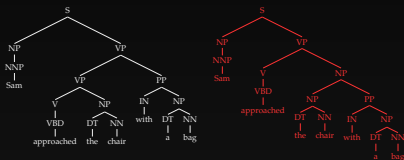
Learning syntax

# Language learning: syntax; in progress

Danny approached the chair with a bag.

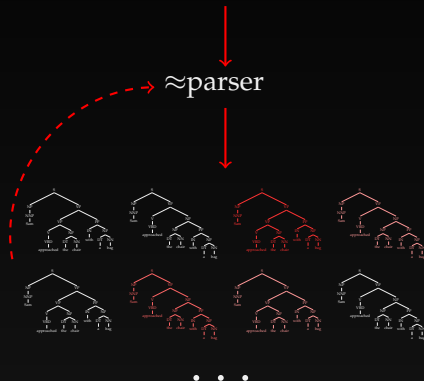


parser



# Language learning: syntax; in progress

Danny approached the chair with a bag.





Courtesy of Elsevier, Inc., <http://www.sciencedirect.com>. Used with permission.  
Source: Pilley, John W., and Alliston K. Reid. "Border collie comprehends object names as verbal referents." *Behavioural processes* 86, no. 2 (2011): 184-195.

Recognition

$$S(\text{sentence}, \text{video})$$

Retrieval

$$\operatorname{argmax}_{v \in V} S(s, v)$$

Generation

$$\operatorname{argmax}_{s \in L} S(s, v)$$

Question answering

$$\operatorname{argmax}_{s \in L} S(Q(s, s_q), v)$$

Disambiguation

$$\operatorname{argmax}_{i \in \text{parser}(s)} S(i, v)$$

Acquiring language

$$\operatorname{argmax}_p \prod_{s, v} S(s(p), v)$$

Images, not videos

$$S(\text{sentence}, \text{video})$$

Translation

Planning

Theory of mind

...



Recognition	$S(\text{sentence}, \text{video})$
Retrieval	$\operatorname{argmax}_{v \in V} S(s, v)$
Generation	$\operatorname{argmax}_{s \in L} S(s, v)$
Question answering	$\operatorname{argmax}_{s \in L} S(Q(s, s_q), v)$
Disambiguation	$\operatorname{argmax}_{i \in \text{parser}(s)} S(i, v)$
Acquiring language	$\operatorname{argmax}_p \prod_{s, v} S(s(p), v)$
Images, not videos	$S(\text{sentence}, \text{Flow}(\text{image}))$
Translation	
Planning	
Theory of mind	
...	

# Single-frame optical flow



© Source Unknown. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>.



© Source Unknown. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>.

# Single-frame optical flow

Input



Flow



Predicted



© Source Unknown. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>.

Recognition	$S(\text{sentence}, \text{video})$
Retrieval	$\operatorname{argmax}_{v \in V} S(s, v)$
Generation	$\operatorname{argmax}_{s \in L} S(s, v)$
Question answering	$\operatorname{argmax}_{s \in L} S(Q(s, s_q), v)$
Disambiguation	$\operatorname{argmax}_{i \in \text{parser}(s)} S(i, v)$
Acquiring language	$\operatorname{argmax}_p \prod_{s, v} S(s(p), v)$
Images, not videos	$S(\text{sentence}, \text{Flow}(\text{image}))$
Translation	$\operatorname{argmin}_{s' \in L_b} f(s, s')$
Planning	
Theory of mind	

...

Recognition	$S(\text{sentence}, \text{video})$
Retrieval	$\operatorname{argmax}_{v \in V} S(s, v)$
Generation	$\operatorname{argmax}_{s \in L} S(s, v)$
Question answering	$\operatorname{argmax}_{s \in L} S(Q(s, s_q), v)$
Disambiguation	$\operatorname{argmax}_{i \in \text{parser}(s)} S(i, v)$
Acquiring language	$\operatorname{argmax}_p \prod_{s, v} S(s(p), v)$
Images, not videos	$S(\text{sentence}, \text{Flow}(\text{image}))$
Translation	$\operatorname{argmin}_{s' \in L_b} \int_v  S(s', v) - S(s, v) $
Planning	
Theory of mind	
...	

# Statistical machine translation

Sam was happy

parallel corpus

Sam a fost fericit<sup>a</sup>

СЭМ БЫЛ<sup>a</sup> СЧАСТЛИВ<sup>a</sup>

In Thai you specify siblings by age not gender.

In English you specify relative time but you don't need to in Chinese.

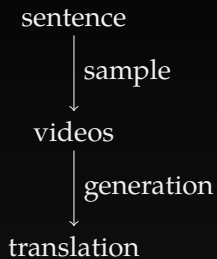
Guugu Yimithirr language only uses absolute directions.

Many languages don't distinguish blue/green.

Swahili specifies color as "the color of X".

In Turkish you have to report if something is hearsay.

# Translation by imagination



Recognition	$S(\text{sentence}, \text{video})$
Retrieval	$\operatorname{argmax}_{v \in V} S(s, v)$
Generation	$\operatorname{argmax}_{s \in L} S(s, v)$
Question answering	$\operatorname{argmax}_{s \in L} S(Q(s, s_q), v)$
Disambiguation	$\operatorname{argmax}_{i \in \text{parser}(s)} S(i, v)$
Acquiring language	$\operatorname{argmax}_p \prod_{s, v} S(s(p), v)$
Images, not videos	$S(\text{sentence}, \text{Flow}(\text{image}))$
Translation	$\operatorname{argmin}_{s' \in L_b} \int_v  S(s', v) - S(s, v) $
Planning	$\operatorname{argmax}_{s \in L} \int_v S(s, v_0 : v : v_n)$
Theory of mind	

...



Recognition	$S(\text{sentence}, \text{video})$
Retrieval	$\operatorname{argmax}_{v \in V} S(s, v)$
Generation	$\operatorname{argmax}_{s \in L} S(s, v)$
Question answering	$\operatorname{argmax}_{s \in L} S(Q(s, s_q), v)$
Disambiguation	$\operatorname{argmax}_{i \in \text{parser}(s)} S(i, v)$
Acquiring language	$\operatorname{argmax}_p \prod_{s, v} S(s(p), v)$
Images, not videos	$S(\text{sentence}, \text{Flow}(\text{image}))$
Translation	$\operatorname{argmin}_{s' \in L_b} \int_v  S(s', v) - S(s, v) $
Planning	$\operatorname{argmax}_{s \in L} \int_v S(s, v_0 : v : v_n)$
Theory of mind	$S(s, v)$

...

Recognition	$S(\text{sentence}, \text{video})$
Retrieval	$\operatorname{argmax}_{v \in V} S(s, v)$
Generation	$\operatorname{argmax}_{s \in L} S(s, v)$
Question answering	$\operatorname{argmax}_{s \in L} S(Q(s, s_q), v)$
Disambiguation	$\operatorname{argmax}_{i \in \text{parser}(s)} S(i, v)$
Acquiring language	$\operatorname{argmax}_p \prod_{s, v} S(s(p), v)$
Images, not videos	$S(\text{sentence}, \text{Flow}(\text{image}))$
Translation	$\operatorname{argmin}_{s' \in L_b} \int_v  S(s', v) - S(s, v) $
Planning	$\operatorname{argmax}_{s \in L} \int_v S(s, v_0 : v : v_n)$
Theory of mind	$S(s, \text{tracks})S(\text{tracks}, v)$

...

Recognition	$S(\text{sentence}, \text{video})$
Retrieval	$\operatorname{argmax}_{v \in V} S(s, v)$
Generation	$\operatorname{argmax}_{s \in L} S(s, v)$
Question answering	$\operatorname{argmax}_{s \in L} S(Q(s, s_q), v)$
Disambiguation	$\operatorname{argmax}_{i \in \text{parser}(s)} S(i, v)$
Acquiring language	$\operatorname{argmax}_p \prod_{s, v} S(s(p), v)$
Images, not videos	$S(\text{sentence}, \text{Flow}(\text{image}))$
Translation	$\operatorname{argmin}_{s' \in L_b} \int_v  S(s', v) - S(s, v) $
Planning	$\operatorname{argmax}_{s \in L} \int_v S(s, v_0 : v : v_n)$
Theory of mind	$S(\text{planner}, \text{tracks})S(s, \text{tracks})S(\text{tracks}, v)$

...

Help	-X Size	+X Size 3	-YX Size	+YX Size 3	Mesh Training	Show Mesh	Clear Image	Quit
Reset	-Z Size	+Z Size 3	-YZ Size	+YZ Size 3	-Mass Threshold	+Mass Threshold 300	Matlab	Record
Test -Theta	+Theta 0.000	Test -Phi	+Phi 28.000	Test -Psi	+psi 0.000	-Center X	+ 0.1326	Start Videos
Test -xi	+xi -5.000	Test -gamma	+gamma 0.000	Test -zeta	+zeta -0.000	-Center Z	+ -0.6103	Stop Videos
-H Span	+H Span 7.9	-V Span	+V Span 1.75	-V Separation	+V Sep 0.9	-Focal Length	+ 1419.615	
-angle	+angle 0	-length	+length 46.700	-height	+height 20.400	Right		
-pan	pan 0.741	-tilt	+tilt 30.136	-bracket	+bracket 0.	-separation	+separation 10.	
Load Ground Truth	Load Result	Load Image	Next	Previous	Show structure	Show Ground Truth	Show Grid	
Get Pose	Get Structure	Cycle Structures	View One	Get Image	Merge Views	Forget Second	Language	Disassemble

741

# The long road ahead . . .

Coherent stories

3D

Forces & contact relations



© sodlvs at Youtube.com. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>.



© Erickson. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>.

# The long road ahead . . .

Coherent stories

3D

Forces & contact relations

Segmentation

Parts and low-level features





© Source Unknown. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>.

# The long road ahead . . .

Coherent stories

3D

Forces & contact relations

Segmentation

Parts and low-level features

Theory of mind

Physics

Modification



© Source Unknown. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>.

# The long road ahead . . .

Coherent stories

3D

Forces & contact relations

Segmentation

Parts and low-level features

Theory of mind

Physics

Modification

The vast majority of verbs: absolve, admire, anger, approve, bark, bend, chase, cheat, complete, concede, discover, fire, follow, fumble, hurry, race, recruit, reject, scratch, steal, taste, want, etc.

Metaphoric extension

*etc.*

# Thanks to many great collaborators

Yevgeni Berzak, Danny Harari, Maximilian Nickel,  
Candance Ross, Victor Carbarera, Santiago Perez,  
Boris Katz, Shimon Ullman, Tomaso Poggio

Siddharth Narayanaswamy, Jeffrey Siskind,  
Sven Dickinson, Song Wang,  
Haonan Yu, Caiming Xiong

Recognition	$S(\text{sentence}, \text{video})$
Retrieval	$\operatorname{argmax}_{v \in V} S(s, v)$
Generation	$\operatorname{argmax}_{s \in L} S(s, v)$
Question answering	$\operatorname{argmax}_{s \in L} S(Q(s, s_q), v)$
Disambiguation	$\operatorname{argmax}_{i \in \text{parser}(s)} S(i, v)$
Acquiring language	$\operatorname{argmax}_p \prod_{s, v} S(s(p), v)$
Images, not videos	$S(\text{sentence}, \text{Flow}(\text{image}))$
Translation	$\operatorname{argmin}_{s' \in L_b} \int_v  S(s', v) - S(s, v) $
Planning	$\operatorname{argmax}_{s \in L} \int_v S(s, v_0 : v : v_n)$
Theory of mind	$S(\text{planner}, \text{tracks})S(s, \text{tracks})S(\text{tracks}, v)$

...

MIT OpenCourseWare  
<https://ocw.mit.edu>

## Resource: Brains, Minds and Machines Summer Course

Tomaso Poggio and Gabriel Kreiman

The following may not correspond to a particular course on MIT OpenCourseWare, but has been provided by the author as an individual learning resource.

For information about citing these materials or our Terms of Use, visit: <https://ocw.mit.edu/terms>.