

# Flexible intelligence



CENTER FOR  
Brains  
Minds +  
Machines

Andrei Barbu



Boris Katz



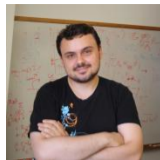
Shimon Ullman



Josh Tenenbaum



Gabriel Kreiman



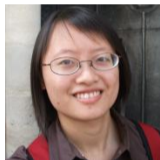
Andrei Barbu



Ignacio Cases



Candace Ross



Yen-Ling Kuo



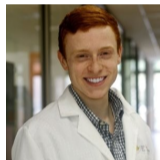
Adam Yaari



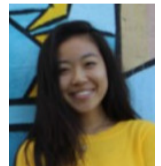
David Mayo



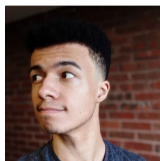
Christopher Wang



Julian Alverio



Emily Cheng



Dylan Sleeper



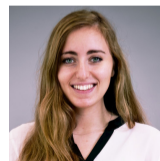
Vighnesh Subramaniam



Aaditya Singh



Ravi Tejwani



Dana Rosenfarb



# Flexible intelligence

# Flexible intelligence

Knew who the adult is talking to

# Flexible intelligence

Knew who the adult is talking to

What “That’s your train” is referring to

# Flexible intelligence

Knew who the adult is talking to

What “That’s your train” is referring to

What trains do on tracks

# Flexible intelligence

Knew who the adult is talking to

What “That’s your train” is referring to

What trains do on tracks

Tried to put trains on tracks



# Flexible intelligence

Knew who the adult is talking to

What “That’s your train” is referring to

What trains do on tracks

Tried to put trains on tracks

Knew what the adult was referring to when they waved their hand

# Flexible intelligence

Knew who the adult is talking to

What “That’s your train” is referring to

What trains do on tracks

Tried to put trains on tracks

Knew what the adult was referring to when they waved their hand

Decided what object to get next

# Flexible intelligence

Knew who the adult is talking to

What “That’s your train” is referring to

What trains do on tracks

Tried to put trains on tracks

Knew what the adult was referring to when they waved their hand

Decided what object to get next

Knew what Santa brought

# Flexible intelligence

Knew who the adult is talking to

What “That’s your train” is referring to

What trains do on tracks

Tried to put trains on tracks

Knew what the adult was referring to when they waved their hand

Decided what object to get next

Knew what Santa brought

Answered questions about objects and abstract concepts

# Flexible intelligence

Knew who the adult is talking to

What “That’s your train” is referring to

What trains do on tracks

Tried to put trains on tracks

Knew what the adult was referring to when they waved their hand

Decided what object to get next

Knew what Santa brought

Answered questions about objects and abstract concepts

Demonstrated “it goes”

# Flexible intelligence

Knew who the adult is talking to

What “That’s your train” is referring to

What trains do on tracks

Tried to put trains on tracks

Knew what the adult was referring to when they waved their hand

Decided what object to get next

Knew what Santa brought

Answered questions about objects and abstract concepts

Demonstrated “it goes”

Understood that “doggie” and “dog” are similar

# Flexible intelligence

Knew who the adult is talking to

What “That’s your train” is referring to

What trains do on tracks

Tried to put trains on tracks

Knew what the adult was referring to when they waved their hand

Decided what object to get next

Knew what Santa brought

Answered questions about objects and abstract concepts

Demonstrated “it goes”

Understood that “doggie” and “dog” are similar

etc.

# Machines on the other hand

*Generated by GPT-2 large*



# Machines on the other hand

*Generated by GPT-2 large*

Q: Leader of Canada A: Trudeau

Q: Leader of Germany A: **Merkel**

**Q: Leader of India A: Modi**

# Machines on the other hand

*Generated by GPT-2 large*

Q: Leader of Canada A: Trudeau

Q: Leader of Germany A: **Merkel**

**Q: Leader of India A: Modi**

Canada is a country that's **become very strong over the last five or six years and has been at the heart of our growth over the last decade or so,"** he told host **Andrew Coyne**

## Machines on the other hand

*Generated by GPT-2 large*

Q: Leader of Canada A: Trudeau

Q: Leader of Germany A: **Merkel**

**Q: Leader of India A: Modi**

Canada is a country that's **become very strong over the last five or six years and has been at the heart of our growth over the last decade or so,"** he told host **Andrew Coyne**

I was on the glideslope and **looked down and there's a lot of wind and clouds and there's lightning coming from the west, so my flight path went off to the left in that direction.**

# Get some things terribly wrong

*GPT-2 large*

The first 3 completions for

I was holding something heavy and my friend . . .

# Get some things terribly wrong

*GPT-2 large*

The first 3 completions for

I was holding something heavy and my friend . . .  
**thought I was getting too close.**

# Get some things terribly wrong

*GPT-2 large*

The first 3 completions for

I was holding something heavy and my friend . . .

**thought I was getting too close.**

**held my back so I wasn't able to grab anything.**

# Get some things terribly wrong

*GPT-2 large*

The first 3 completions for

I was holding something heavy and my friend . . .

**thought I was getting too close.**

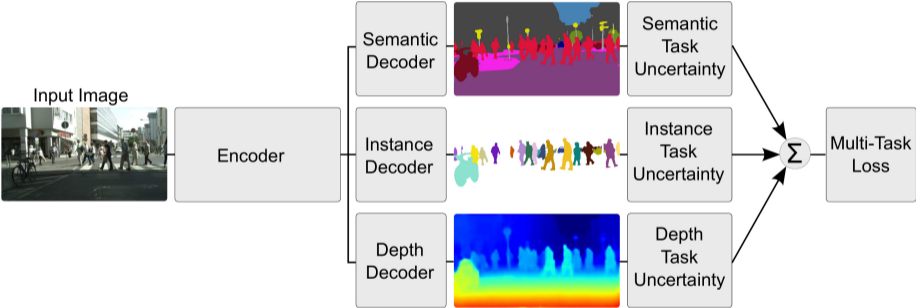
**held my back so I wasn't able to grab anything.**

**took it and hit me in the head,"**

# Multi-task Vision



# Multi-task Vision







Recognition

Retrieval

Recognition

Retrieval

Generation

Recognition

Retrieval

Generation

Question answering

Recognition

Retrieval

Generation

Question answering

Disambiguation

Recognition

Retrieval

Generation

Question answering

Disambiguation

Language acquisition



Recognition

Retrieval

Generation

Question answering

Disambiguation

Language acquisition

Paraphrasing

Recognition

Retrieval

Generation

Question answering

Disambiguation

Language acquisition

Paraphrasing

Translation

Recognition

Retrieval

Generation

Question answering

Disambiguation

Language acquisition

Paraphrasing

Translation

Common sense reasoning

Recognition

Retrieval

Generation

Question answering

Disambiguation

Language acquisition

Paraphrasing

Translation

Common sense reasoning

Planning

Recognition

Retrieval

Generation

Question answering

Disambiguation

Language acquisition

Paraphrasing

Translation

Common sense reasoning

Planning

Command following

Recognition

Retrieval

Generation

Question answering

Disambiguation

Language acquisition

Paraphrasing

Translation

Common sense reasoning

Planning

Command following

...

Recognition

Retrieval

Generation

Question answering

Disambiguation

Language acquisition

Paraphrasing

Translation

Common sense reasoning

Planning

Command following

...

Computer vision

NLP

Robotics

AI

# Representation

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



# Representation

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

Sentence  $\rightarrow$  First order temporal logic  $\rightarrow$  video detector

# Representation

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

Sentence  $\rightarrow$  First order temporal logic  $\rightarrow$  video detector

Danny and Andrei move a chair.

# Representation

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

Sentence  $\rightarrow$  First order temporal logic  $\rightarrow$  video detector

Danny and Andrei move a chair.

$\exists xyz \text{ chair}(x), \text{person}(y), \text{person}(z), y \neq z, \text{move}(y, x), \text{move}(z, x)$

# Representation

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

Sentence  $\rightarrow$  First order temporal logic  $\rightarrow$  video detector

Danny and Andrei move a chair.

$\exists xyz \text{ chair}(x), \text{person}(y), \text{person}(z), y \neq z, \text{move}(y, x), \text{move}(z, x)$

Don't build in attributes like *shape*, *color*, etc.

# Representation

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

Sentence  $\rightarrow$  First order temporal logic  $\rightarrow$  video detector

Danny and Andrei move a chair.

$\exists xyz \text{ chair}(x), \text{person}(y), \text{person}(z), y \neq z, \text{move}(y, x), \text{move}(z, x)$

Don't build in attributes like *shape*, *color*, etc.

Build in mechanics! Physics, Interactions, ToM, out of domain goals, etc.

# Representation

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

Sentence  $\rightarrow$  First order temporal logic  $\rightarrow$  video detector

Danny and Andrei move a chair.

$\exists xyz \text{ chair}(x), \text{person}(y), \text{person}(z), y \neq z, \text{move}(y, x), \text{move}(z, x)$

Don't build in attributes like *shape*, *color*, etc.

Build in mechanics! Physics, Interactions, ToM, out of domain goals, etc.

Attributes exist because they matter, mechanisms are built in

# Representation

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

Sentence  $\rightarrow$  First order temporal logic  $\rightarrow$  video detector

Danny and Andrei move a chair.

$\exists xyz \text{ chair}(x), \text{person}(y), \text{person}(z), y \neq z, \text{move}(y, x), \text{move}(z, x)$

Don't build in attributes like *shape*, *color*, etc.

Build in mechanics! Physics, Interactions, ToM, out of domain goals, etc.

Attributes exist because they matter, mechanisms are built in

Compositions of networks or graphical models:

# Representation

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

Sentence  $\rightarrow$  First order temporal logic  $\rightarrow$  video detector

Danny and Andrei move a chair.

$\exists xyz \text{ chair}(x), \text{person}(y), \text{person}(z), y \neq z, \text{move}(y, x), \text{move}(z, x)$

Don't build in attributes like *shape*, *color*, etc.

Build in mechanics! Physics, Interactions, ToM, out of domain goals, etc.

Attributes exist because they matter, mechanisms are built in

Compositions of networks or graphical models:

...

Yu, Siddharth, Barbu, Siskind JAIR 2015



# Representation

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

Sentence  $\rightarrow$  First order temporal logic  $\rightarrow$  video detector

Danny and Andrei move a chair.

$$\exists xyz \text{ chair}(x), \text{ person}(y), \text{ person}(z), y \neq z, \text{ move}(y, x), \text{ move}(z, x)$$

Don't build in attributes like *shape*, *color*, etc.

Build in mechanics! Physics, Interactions, ToM, out of domain goals, etc.

Attributes exist because they matter, mechanisms are built in

Compositions of networks or graphical models:

x-tracker

y-tracker

z-tracker



...

Yu, Siddharth, Barbu, Siskind JAIR 2015

# Representation

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

Sentence  $\rightarrow$  First order temporal logic  $\rightarrow$  video detector

Danny and Andrei move a chair.

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

Don't build in attributes like *shape*, *color*, etc.

Build in mechanics! Physics, Interactions, ToM, out of domain goals, etc.

Attributes exist because they matter, mechanisms are built in

Compositions of networks or graphical models:

chair person person  $\neq$  move move

x-tracker

y-tracker

z-tracker



...

Yu, Siddharth, Barbu, Siskind JAIR 2015

# Representation

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

Sentence  $\rightarrow$  First order temporal logic  $\rightarrow$  video detector

Danny and Andrei move a chair.

$$\exists xyz \text{ chair}(x), \text{ person}(y), \text{ person}(z), y \neq z, \text{ move}(y, x), \text{ move}(z, x)$$

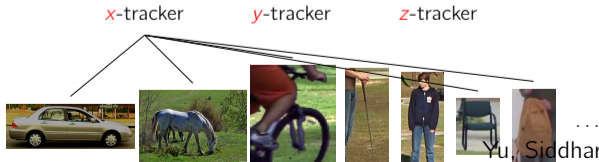
Don't build in attributes like *shape*, *color*, etc.

Build in mechanics! Physics, Interactions, ToM, out of domain goals, etc.

Attributes exist because they matter, mechanisms are built in

Compositions of networks or graphical models:

chair person person  $\neq$  move move



# Representation

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

Sentence  $\rightarrow$  First order temporal logic  $\rightarrow$  video detector

Danny and Andrei move a chair.

$$\exists xyz \text{ chair}(x), \text{ person}(y), \text{ person}(z), y \neq z, \text{ move}(y, x), \text{ move}(z, x)$$

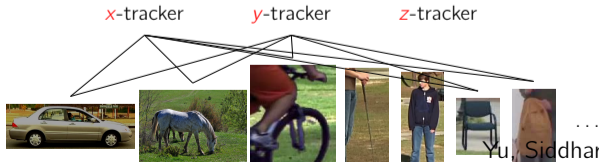
Don't build in attributes like *shape*, *color*, etc.

Build in mechanics! Physics, Interactions, ToM, out of domain goals, etc.

Attributes exist because they matter, mechanisms are built in

Compositions of networks or graphical models:

chair person person  $\neq$  move move



# Representation

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

Sentence  $\rightarrow$  First order temporal logic  $\rightarrow$  video detector

Danny and Andrei move a chair.

$$\exists xyz \text{ chair}(x), \text{ person}(y), \text{ person}(z), y \neq z, \text{ move}(y, x), \text{ move}(z, x)$$

Don't build in attributes like *shape*, *color*, etc.

Build in mechanics! Physics, Interactions, ToM, out of domain goals, etc.

Attributes exist because they matter, mechanisms are built in

Compositions of networks or graphical models:

chair person person  $\neq$  move move



# Representation

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

Sentence  $\rightarrow$  First order temporal logic  $\rightarrow$  video detector

Danny and Andrei move a chair.

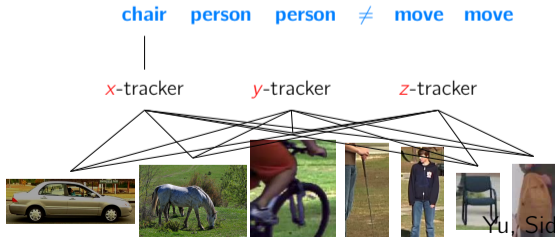
$$\exists xyz \text{ chair}(x), \text{ person}(y), \text{ person}(z), y \neq z, \text{ move}(y, x), \text{ move}(z, x)$$

Don't build in attributes like *shape*, *color*, etc.

Build in mechanics! Physics, Interactions, ToM, out of domain goals, etc.

Attributes exist because they matter, mechanisms are built in

Compositions of networks or graphical models:



# Representation

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

Sentence  $\rightarrow$  First order temporal logic  $\rightarrow$  video detector

Danny and Andrei move a chair.

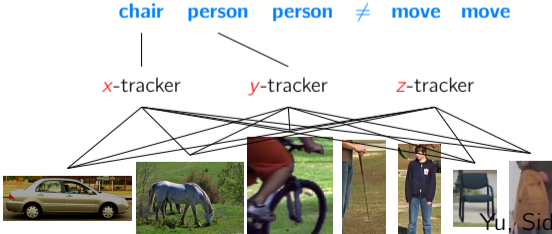
$$\exists xyz \text{ chair}(x), \text{ person}(y), \text{ person}(z), y \neq z, \text{ move}(y, x), \text{ move}(z, x)$$

Don't build in attributes like *shape*, *color*, etc.

Build in mechanics! Physics, Interactions, ToM, out of domain goals, etc.

Attributes exist because they matter, mechanisms are built in

Compositions of networks or graphical models:



Yu, Siddharth, Barbu, Siskind JAIR 2015

# Representation

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

Sentence  $\rightarrow$  First order temporal logic  $\rightarrow$  video detector

Danny and Andrei move a chair.

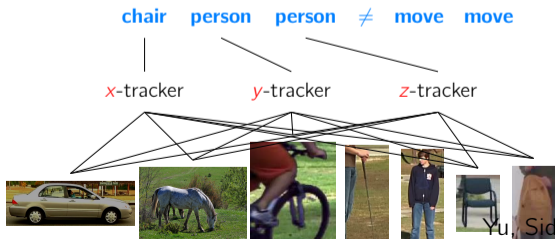
$$\exists xyz \text{ chair}(x), \text{ person}(y), \text{ person}(z), y \neq z, \text{ move}(y, x), \text{ move}(z, x)$$

Don't build in attributes like *shape*, *color*, etc.

Build in mechanics! Physics, Interactions, ToM, out of domain goals, etc.

Attributes exist because they matter, mechanisms are built in

Compositions of networks or graphical models:





# Representation

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

Sentence  $\rightarrow$  First order temporal logic  $\rightarrow$  video detector

Danny and Andrei move a chair.

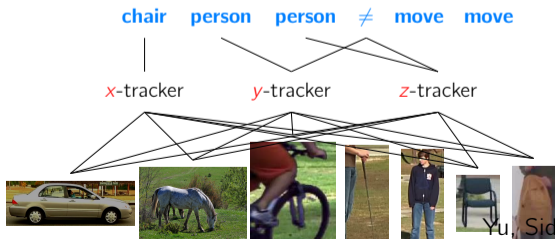
$\exists xyz \text{ chair}(x), \text{person}(y), \text{person}(z), y \neq z, \text{move}(y, x), \text{move}(z, x)$

Don't build in attributes like *shape*, *color*, etc.

Build in mechanics! Physics, Interactions, ToM, out of domain goals, etc.

Attributes exist because they matter, mechanisms are built in

Compositions of networks or graphical models:



# Representation

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

Sentence  $\rightarrow$  First order temporal logic  $\rightarrow$  video detector

Danny and Andrei move a chair.

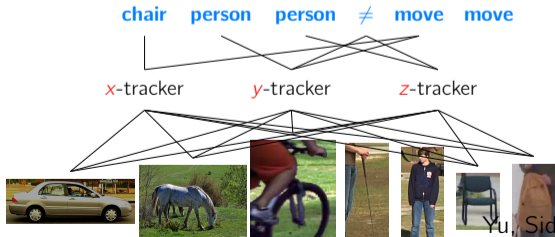
$$\exists xyz \text{ chair}(x), \text{ person}(y), \text{ person}(z), y \neq z, \text{ move}(y, x), \text{ move}(z, x)$$

Don't build in attributes like *shape*, *color*, etc.

Build in mechanics! Physics, Interactions, ToM, out of domain goals, etc.

Attributes exist because they matter, mechanisms are built in

Compositions of networks or graphical models:



Yu, Siddharth, Barbu, Siskind JAIR 2015

# Representation

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

Sentence  $\rightarrow$  First order temporal logic  $\rightarrow$  video detector

Danny and Andrei move a chair.

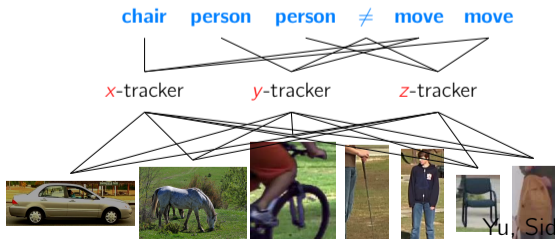
$$\exists xyz \text{ chair}(x), \text{ person}(y), \text{ person}(z), y \neq z, \text{ move}(y, x), \text{ move}(z, x)$$

Don't build in attributes like *shape*, *color*, etc.

Build in mechanics! Physics, Interactions, ToM, out of domain goals, etc.

Attributes exist because they matter, mechanisms are built in

Compositions of networks or graphical models:



Retrieval

Generation

Question answering

Disambiguation

Language acquisition

Paraphrasing

Translation

Common sense reasoning

Planning

Command following

...

Recognition

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

Narayanaswamy *et al.* 2014

Retrieval

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

Barret *et al.* 2016

Generation

Question answering

Disambiguation

Language acquisition

Paraphrasing

Translation

Common sense reasoning

Planning

Command following

...

# Sentential retrieval

# Sentential retrieval



Recognition

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

Narayanaswamy *et al.* 2014

Retrieval

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

Barret *et al.* 2016

Generation

Question answering

Disambiguation

Language acquisition

Paraphrasing

Translation

Common sense reasoning

Planning

Command following

...



Recognition

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

Narayanaswamy *et al.* 2014

Retrieval

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

Barret *et al.* 2016

Generation

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

Yu *et al.* 2015, N. *et al.* 2014

Question answering

Disambiguation

Language acquisition

Paraphrasing

Translation

Common sense reasoning

Planning

Command following

...

# Generating sentences

# Generating sentences

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

# 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

# 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

∅

# 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

“carried”

# 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”

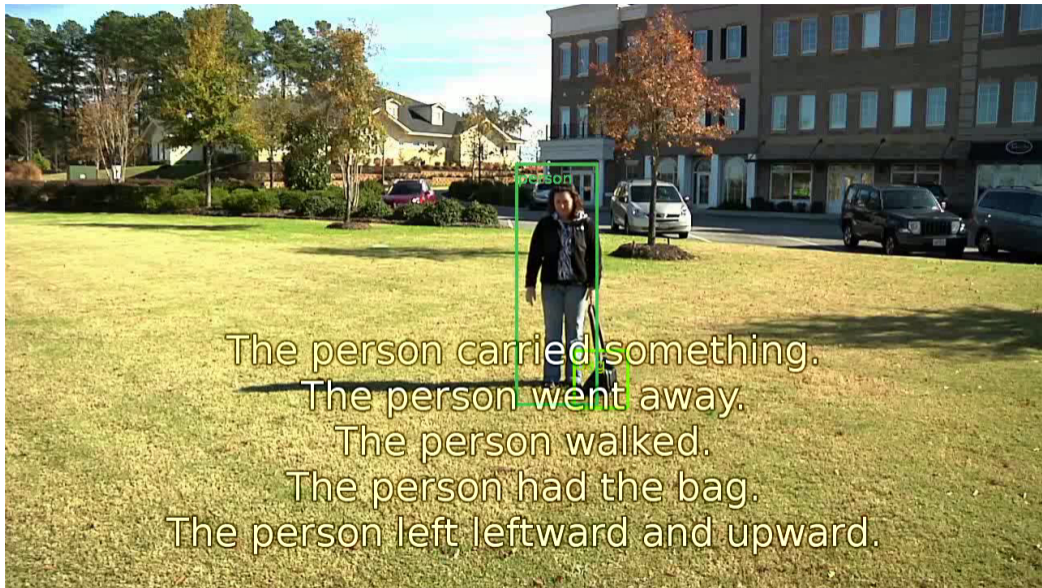
# Generating sentences

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

147,123,874,800 sentences without recursion

“the person carried the backpack”





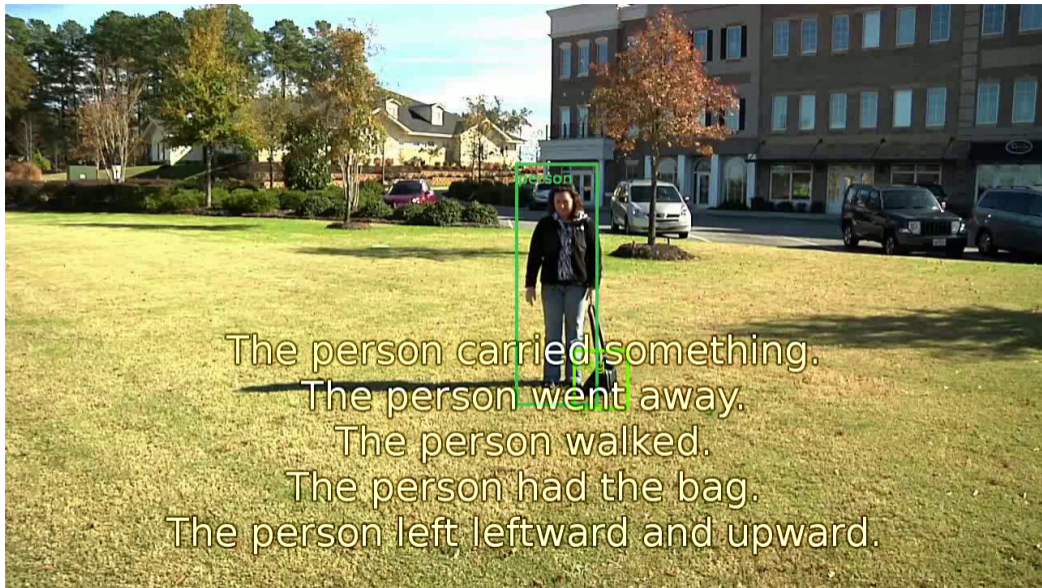
The person carried something.

The person went away.

The person walked.

The person had the bag.

The person left leftward and upward.



The person carried something.

The person went away.

The person walked.

The person had the bag.

The person left leftward and upward.

Recognition

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

Narayanaswamy *et al.* 2014

Retrieval

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

Barret *et al.* 2016

Generation

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

Yu *et al.* 2015, N. *et al.* 2014

Question answering

Disambiguation

Language acquisition

Paraphrasing

Translation

Common sense reasoning

Planning

Command following

...

Recognition

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

Narayanaswamy *et al.* 2014

Retrieval

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

Barret *et al.* 2016

Generation

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

Yu *et al.* 2015, N. *et al.* 2014

Question answering

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

Barbu *et al.* in prep.

Disambiguation

Language acquisition

Paraphrasing

Translation

Common sense reasoning

Planning

Command following

...

# Question answering

## Question answering



## Question answering



## Question answering



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



## Question answering



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

## 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.



I saw the man with the telescope.

I saw the man with the telescope.



I saw the man with the telescope.





Danny looked at the man with a telescope.




Danny looked at the man with a telescope.



Danny has the telescope

Danny looked at the man with a telescope.



Danny has the telescope

The man has the telescope

Danny looked at the man with a telescope.

Danny has the telescope

The man has the telescope



Danny looked at the man with a telescope.

Danny has the telescope

The man has the telescope



Danny looked at the man with a telescope.

Danny has the telescope

The man has the telescope



# Ambiguities

# Ambiguities

PP Attachment

Danny looked at the man with a telescope.



# Ambiguities

PP Attachment

Andrei approached the person holding a green chair.

VP Attachment





# Ambiguities

PP Attachment

VP Attachment

Conjunction

Danny and Andrei picked up the yellow bag and chair.



# Ambiguities

PP Attachment

VP Attachment

Conjunction

Logical Form

Someone put down the bags.



# Ambiguities

PP Attachment

VP Attachment

Conjunction

Logical Form

Anaphora

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



# Ambiguities

PP Attachment

VP Attachment

Conjunction

Logical Form

Anaphora

Ellipsis

Danny left Andrei. Also Yevgeni.



Recognition	$P(\text{sentence}, \text{video})$	Narayanaswamy <i>et al.</i> 2014
Retrieval	$\operatorname{argmax}_{v \in V} P(s, v)$	Barret <i>et al.</i> 2016
Generation	$\operatorname{argmax}_{s \in L} P(s, v)$	Yu <i>et al.</i> 2015, N. <i>et al.</i> 2014
Question answering	$\operatorname{argmax}_{s \in L} P(Q(s, s_q), v)$	Barbu <i>et al.</i> in prep.
Disambiguation	$\operatorname{argmax}_{i \in \text{parser}(s)} P(i, v)$	Berzak <i>et al.</i> 2015
Language acquisition		
Paraphrasing		
Translation		
Common sense reasoning		
Planning		
Command following		

Recognition	$P(\text{sentence, video})$	Narayanaswamy <i>et al.</i> 2014
Retrieval	$\operatorname{argmax}_{v \in V} P(s, v)$	Barret <i>et al.</i> 2016
Generation	$\operatorname{argmax}_{s \in L} P(s, v)$	Yu <i>et al.</i> 2015, N. <i>et al.</i> 2014
Question answering	$\operatorname{argmax}_{s \in L} P(Q(s, s_q), v)$	Barbu <i>et al.</i> in prep.
Disambiguation	$\operatorname{argmax}_{i \in \text{parser}(s)} P(i, v)$	Berzak <i>et al.</i> 2015
Language acquisition	$\operatorname{argmax}_{\theta} \prod_{s, v} P(s(\theta), v)$	Yu <i>et al.</i> 2015, Ross <i>et al.</i> 2018
Paraphrasing		
Translation		
Common sense reasoning		
Planning		
Command following		

# Language acquisition from 10,000 feet

# Language acquisition from 10,000 feet

*Take this apple.*



# Language acquisition from 10,000 feet

*Take this apple.*



# Language acquisition from 10,000 feet

( *Take this apple.* ,



)

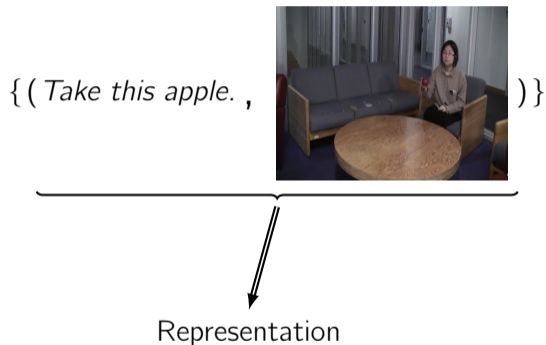
# Language acquisition from 10,000 feet

{ (*Take this apple.* ,

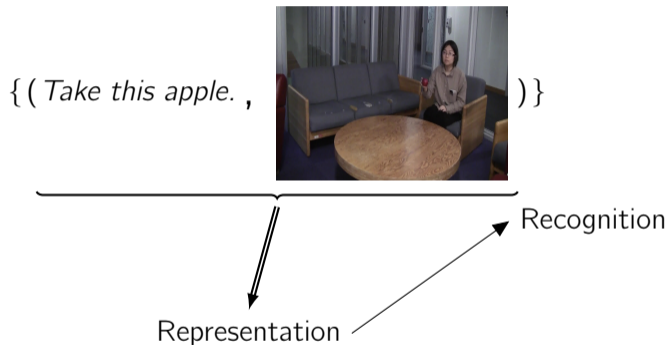


)}

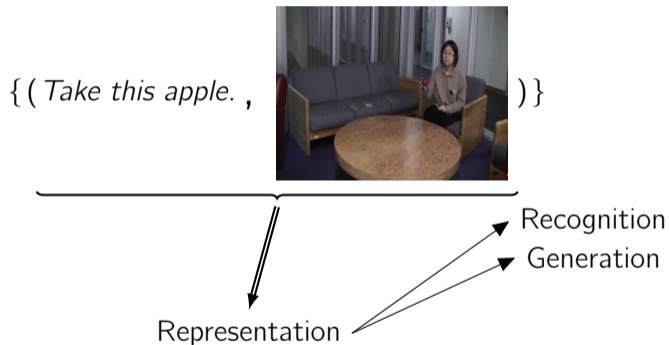
# Language acquisition from 10,000 feet



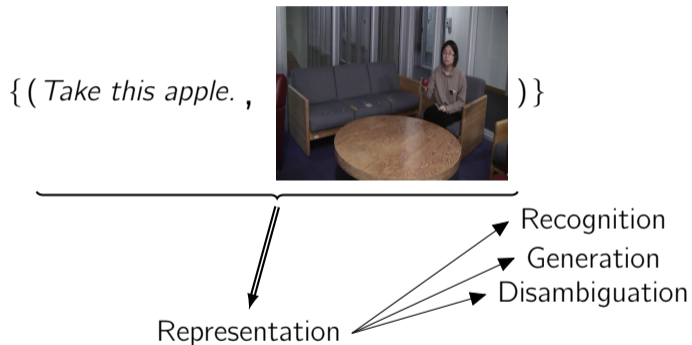
# Language acquisition from 10,000 feet



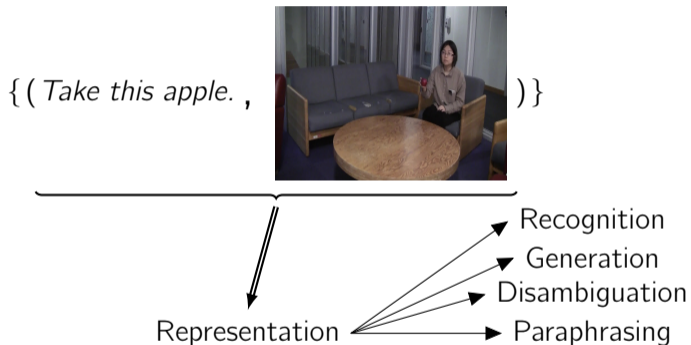
# Language acquisition from 10,000 feet



# Language acquisition from 10,000 feet

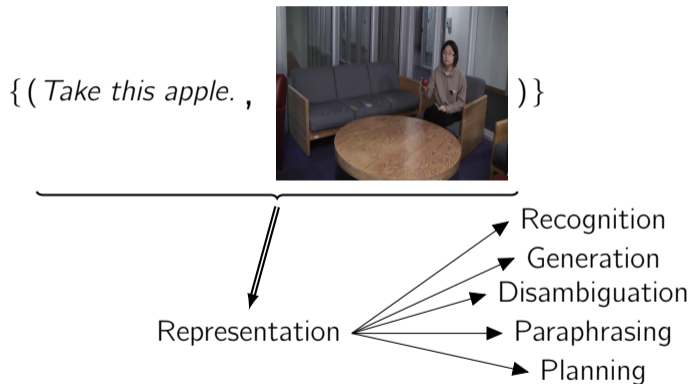


# Language acquisition from 10,000 feet

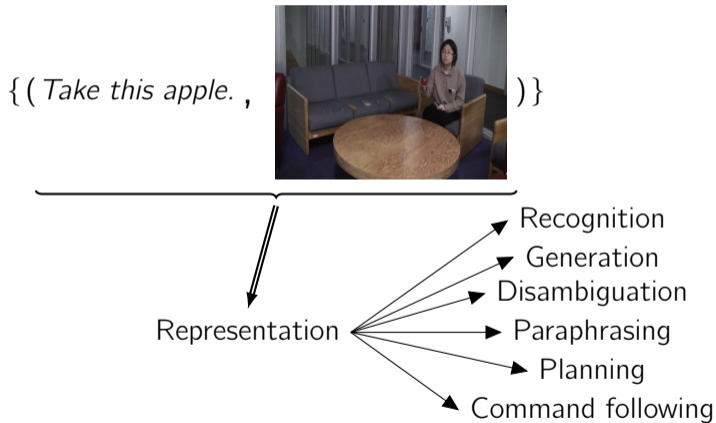




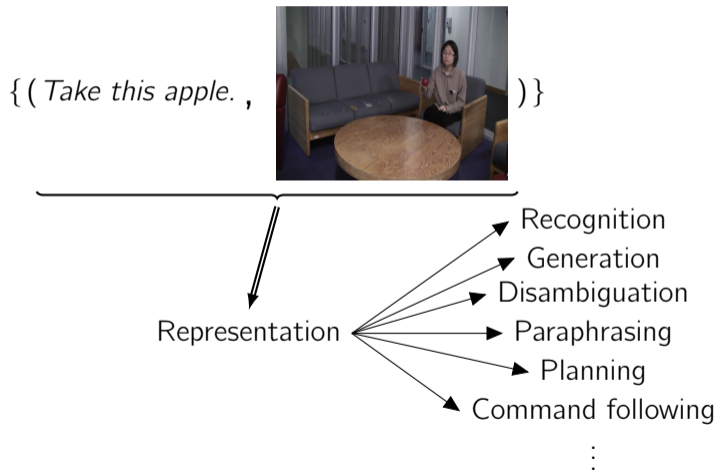
# Language acquisition from 10,000 feet



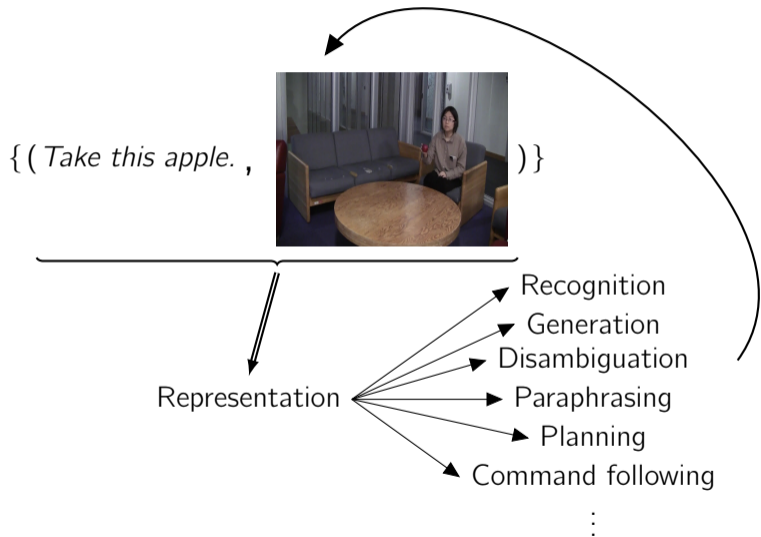
# Language acquisition from 10,000 feet



# Language acquisition from 10,000 feet

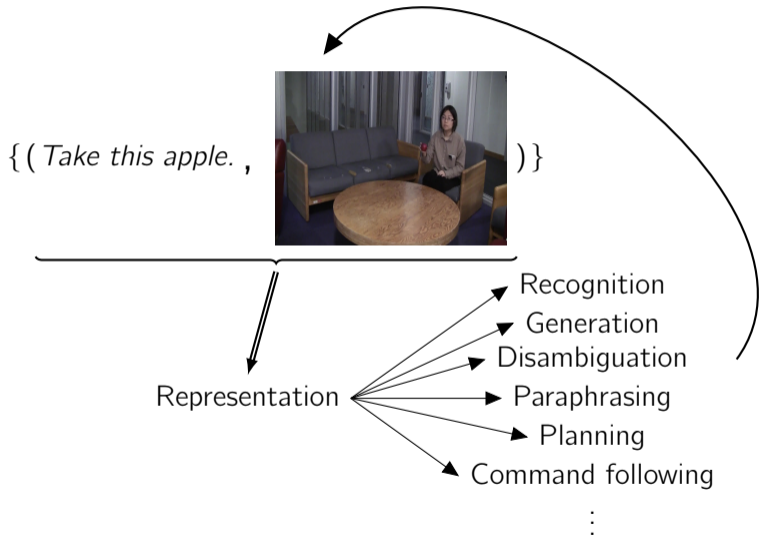


# Language acquisition from 10,000 feet



# Language acquisition from 10,000 feet

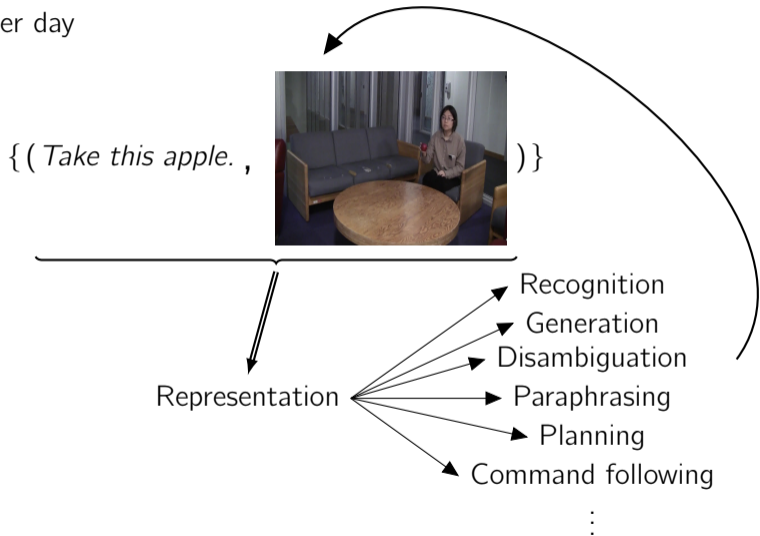
By age 4



# Language acquisition from 10,000 feet

By age 4

heard  $\approx$ 12000 words per day

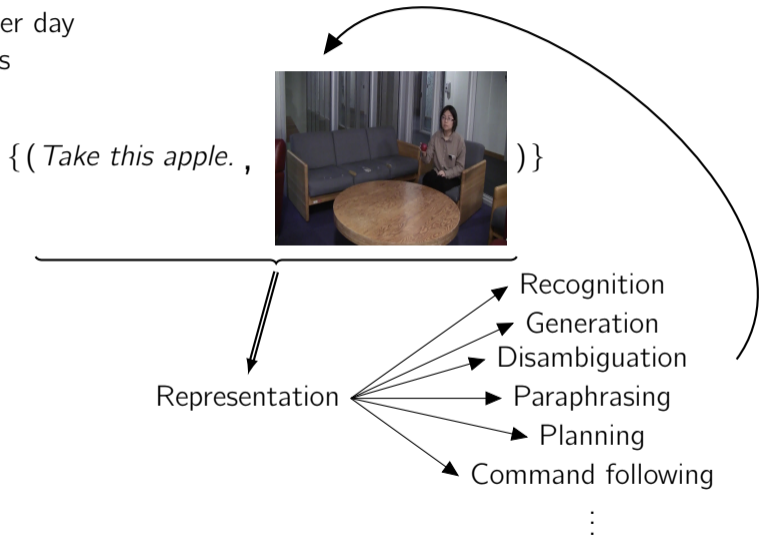


# Language acquisition from 10,000 feet

By age 4

heard  $\approx 12000$  words per day

$\approx 1$ -2 million utterances



# Language acquisition from 10,000 feet

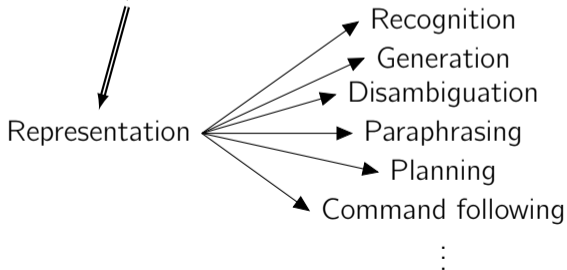
By age 4

heard  $\approx 12000$  words per day

$\approx 1$ -2 million utterances

$\approx 550$  daily conversational turns

{ (*Take this apple.* , ) }





# Grounded language acquisition

# Grounded language acquisition

Danny approached the chair with a bag.

# Grounded language acquisition

Danny approached the chair with a bag.

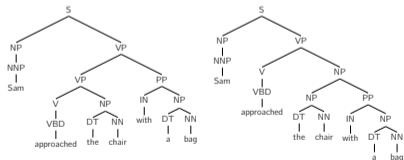


parser

# Grounded language acquisition

Danny approached the chair with a bag.

parser

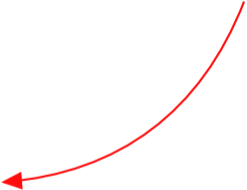
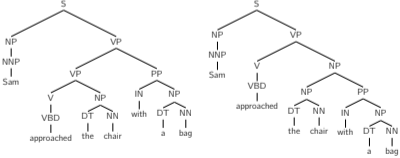


# Grounded language acquisition

Danny approached the chair with a bag.



parser

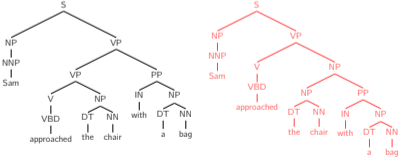


# Grounded language acquisition

Danny approached the chair with a bag.



parser

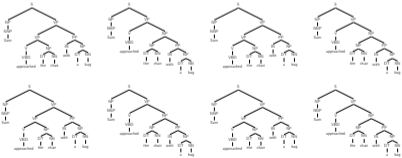


# Grounded language acquisition

Danny approached the chair with a bag.



parser



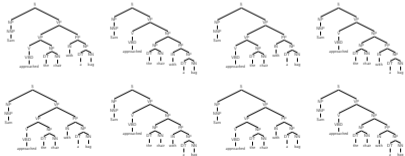
• • •

# Grounded language acquisition

Danny approached the chair with a bag.



≈ parser



• • •



# Grounded language acquisition

Danny approached the chair with a bag.



≈ parser



# Grounded language acquisition

Danny approached the chair with a bag.



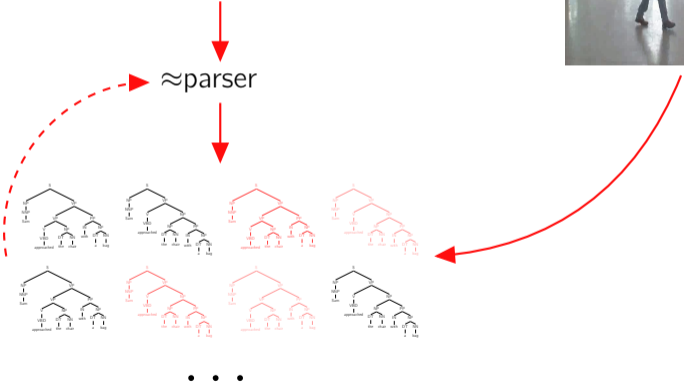
≈ parser



• • •

# Grounded language acquisition

Danny approached the chair with a bag.



# The acquired syntactic & semantic parser

A CCG-based parser with an acquired lexicon and  
a small network that ranks derivations

# The acquired syntactic & semantic parser

A CCG-based parser with an acquired lexicon and  
a small network that ranks derivations

*She places the toy car down on the table.*

# The acquired syntactic & semantic parser

A CCG-based parser with an acquired lexicon and  
a small network that ranks derivations

*She places the toy car down on the table.*

$\lambda xyz.person  $x$ , put-down  $x$   $y$ , toy  $y$ , car  $y$ , table  $z$ , on  $y$   $z$$

# The acquired syntactic & semantic parser

A CCG-based parser with an acquired lexicon and  
a small network that ranks derivations

*She places the toy car down on the table.*

$\lambda xyz.person  $x$ , put-down  $x$   $y$ , toy  $y$ , car  $y$ , table  $z$ , on  $y$   $z$$

Fully-supervised: 93% accuracy

Unsupervised:  $\approx 1\%$  accuracy

Ours, videos without annotations: 60% accuracy

# Learning a parser-generator pair using a robotic simulator

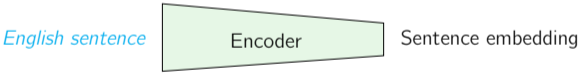


# Learning a parser-generator pair using a robotic simulator

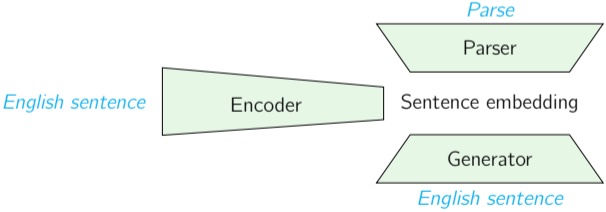


*English sentence*

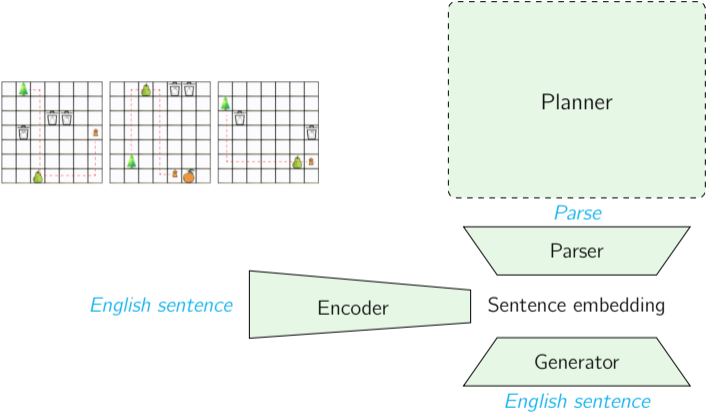
# Learning a parser-generator pair using a robotic simulator



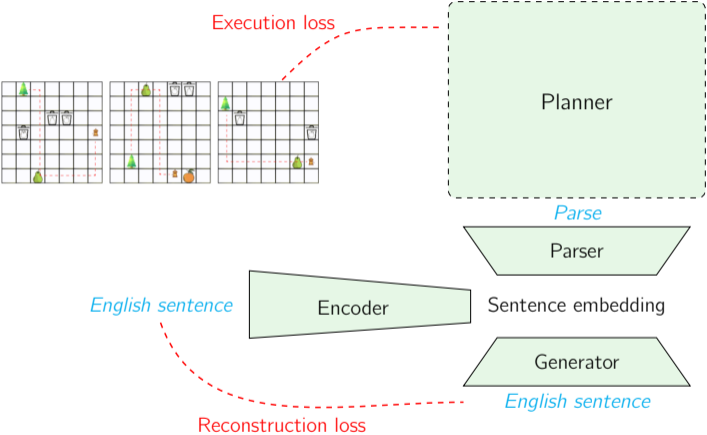
# Learning a parser-generator pair using a robotic simulator



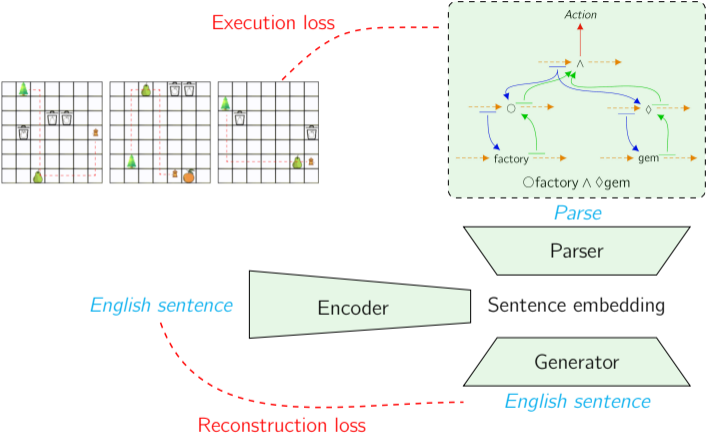
# Learning a parser-generator pair using a robotic simulator



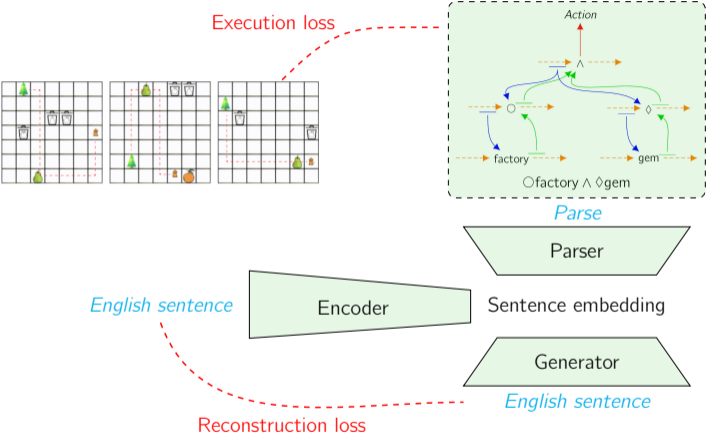
# Learning a parser-generator pair using a robotic simulator



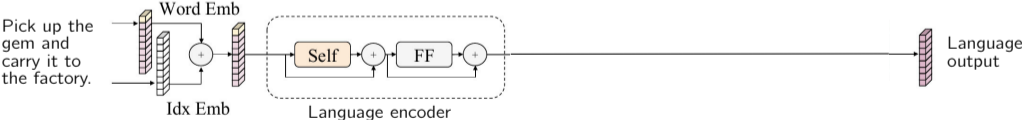
# Learning a parser-generator pair using a robotic simulator



# Use a robot to discover the structure of language

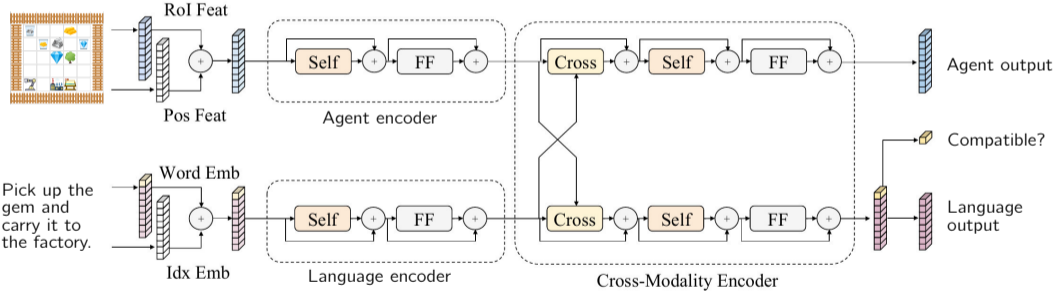


# Language models that interact physically

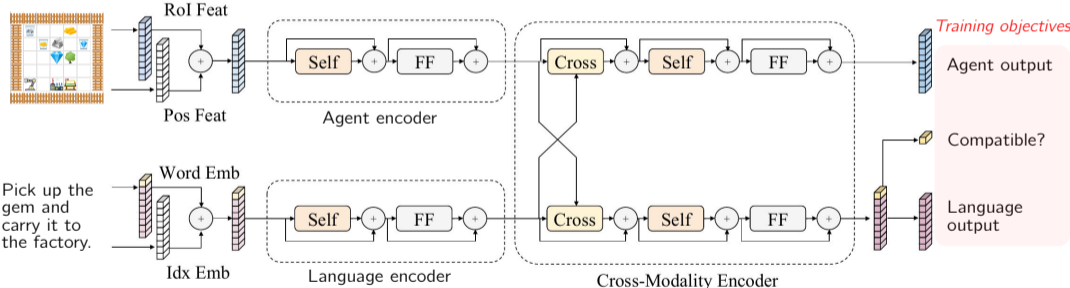




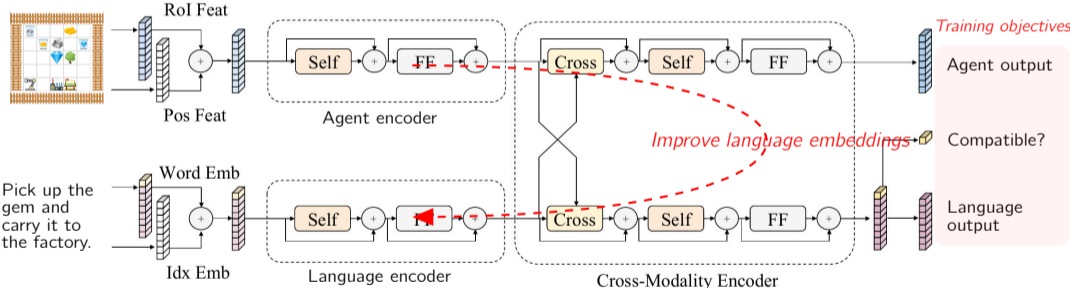
# Language models that interact physically



# Language models that interact physically



# Language models that interact physically



# What about social interactions?

We **have no** social simulators

# What about social interactions?

We **now have** social simulators

PHASE: PHysically-grounded Abstract Social Events for Machine Social Perception

# What about social interactions?

We **now have** social simulators

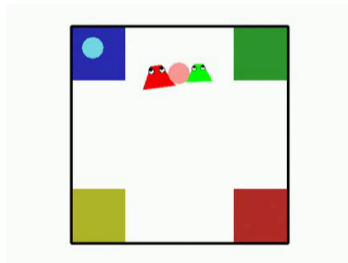
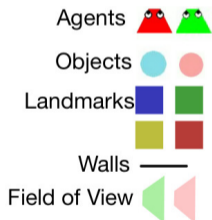
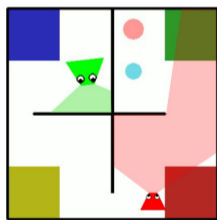
PHASE: PHysically-grounded Abstract Social Events for Machine Social Perception



# What about social interactions?

We **now have** social simulators

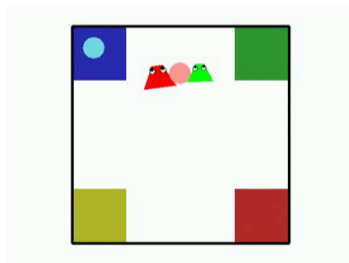
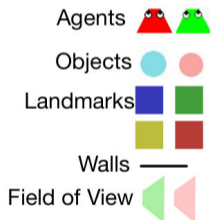
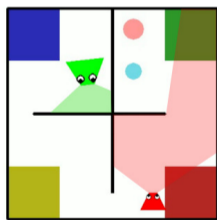
PHASE: PHysically-grounded Abstract Social Events for Machine Social Perception



# What about social interactions?

We **now have** social simulators

PHASE: PHysically-grounded Abstract Social Events for Machine Social Perception



We have no mathematical theories theories of social interactions





Pilley and Reid 2011



Pilley and Reid 2011



Pilley and Reid 2011

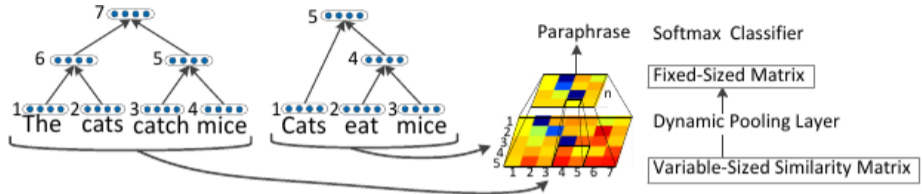
Recognition	$P(\text{sentence, video})$	Narayanaswamy <i>et al.</i> 2014
Retrieval	$\operatorname{argmax}_{v \in V} P(s, v)$	Barret <i>et al.</i> 2016
Generation	$\operatorname{argmax}_{s \in L} P(s, v)$	Yu <i>et al.</i> 2015, N. <i>et al.</i> 2014
Question answering	$\operatorname{argmax}_{s \in L} P(Q(s, s_q), v)$	Barbu <i>et al.</i> in prep.
Disambiguation	$\operatorname{argmax}_{i \in \text{parser}(s)} P(i, v)$	Berzak <i>et al.</i> 2015
Language acquisition	$\operatorname{argmax}_{\theta} \prod_{s, v} P(s(\theta), v)$	Yu <i>et al.</i> 2015, Ross <i>et al.</i> 2018
Paraphrasing		
Translation		
Common sense reasoning		
Planning		
Command following		

Recognition	$P(\text{sentence}, \text{video})$	Narayanaswamy <i>et al.</i> 2014
Retrieval	$\operatorname{argmax}_{v \in V} P(s, v)$	Barret <i>et al.</i> 2016
Generation	$\operatorname{argmax}_{s \in L} P(s, v)$	Yu <i>et al.</i> 2015, N. <i>et al.</i> 2014
Question answering	$\operatorname{argmax}_{s \in L} P(Q(s, s_q), v)$	Barbu <i>et al.</i> in prep.
Disambiguation	$\operatorname{argmax}_{i \in \text{parser}(s)} P(i, v)$	Berzak <i>et al.</i> 2015
Language acquisition	$\operatorname{argmax}_{\theta} \prod_{s, v} P(s(\theta), v)$	Yu <i>et al.</i> 2015, Ross <i>et al.</i> 2018
Paraphrasing	$\int_v  P(s, v) - P(s', v) $	Mao <i>et al.</i> in review
Translation		
Common sense reasoning		
Planning		
Command following		

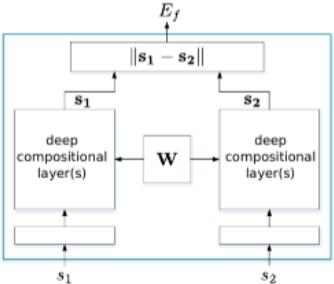
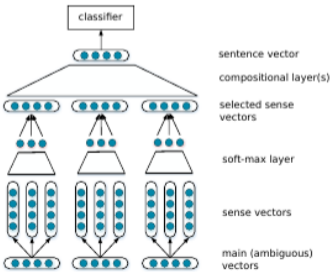
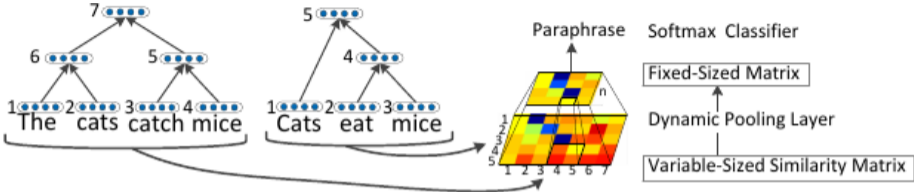
# Paraphrasing today

Socher *et al.* 2011, Cheng and Kartsaklis 2015

# Paraphrasing today



# Paraphrasing today



Socher *et al.* 2011, Cheng and Kartsaklis 2015



# Paraphrasing with vision

# Paraphrasing with vision

$$s \stackrel{?}{\Rightarrow} s'$$

# Paraphrasing with vision

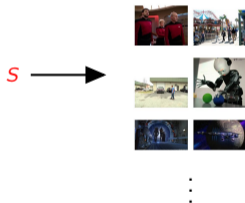
$$s \stackrel{?}{\Rightarrow} s'$$

$s$

# Paraphrasing with vision

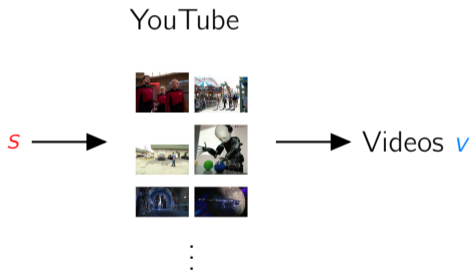
$$s \xrightarrow{?} s'$$

YouTube



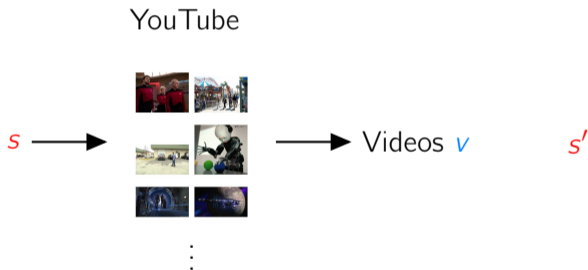
# Paraphrasing with vision

$$s \stackrel{?}{\Rightarrow} s'$$



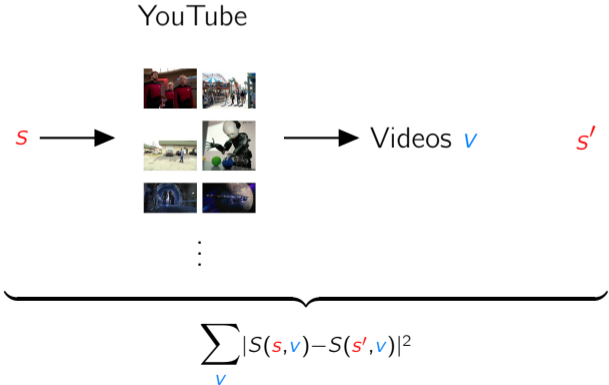
# Paraphrasing with vision

$$s \stackrel{?}{\Rightarrow} s'$$



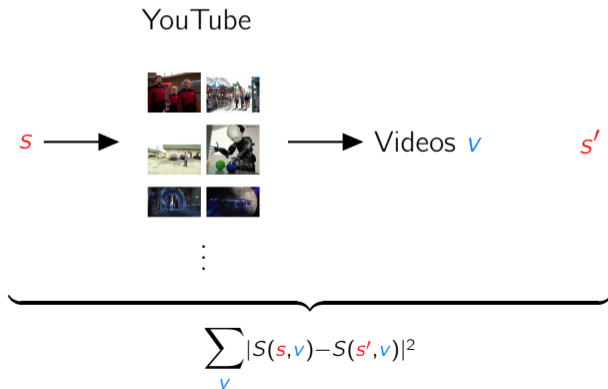
# Paraphrasing with vision

$$s \stackrel{?}{\Rightarrow} s'$$



# Paraphrasing with vision

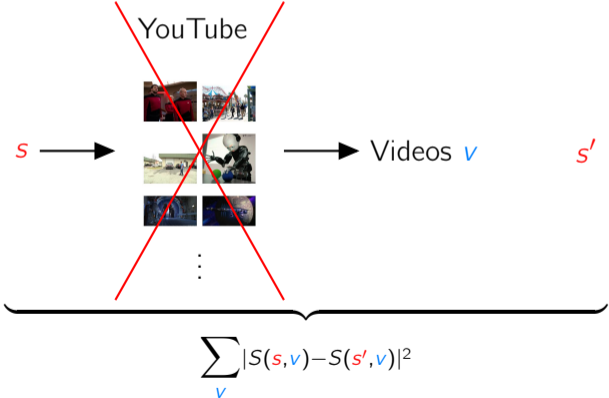
$$s \stackrel{?}{\Rightarrow} s'$$





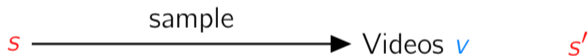
# Paraphrasing with vision

$$s \stackrel{?}{\Rightarrow} s'$$



# Paraphrasing with vision

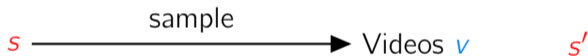
$$s \stackrel{?}{\Rightarrow} s'$$



$$\int_v |S(s, v) - S(s', v)|^2$$

# Paraphrasing with imagination

$$s \stackrel{?}{\Rightarrow} s'$$



$$\int_v |S(s,v) - S(s',v)|^2$$

A large horizontal brace is positioned above the integral symbol, extending from the left side of the diagram to the right side, encompassing the 'sample' process and the 'Videos v' part.

## Generated “videos”

Alice carried the chair away from the backpack.



Mao, Katz, Barbu; IJCAI in review

## Generated “videos”

Alice carried the chair away from the backpack.



Mao, Katz, Barbu; IJCAI in review

## Comparing sentences vs Parikh *et al.* 2016 + ELMo

Does the sentence above imply the one below, and vice versa?

# Comparing sentences vs Parikh *et al.* 2016 + ELMo

Does the sentence above imply the one below, and vice versa?

Ground Ours Theirs

# Comparing sentences vs Parikh *et al.* 2016 + ELMo

Does the sentence above imply the one below, and vice versa?

Ground Ours Theirs

*Alice carried the chair*

*Alice held the chair*



# Comparing sentences vs Parikh *et al.* 2016 + ELMo

Does the sentence above imply the one below, and vice versa?

*Alice carried the chair*  
*Alice held the chair*

Ground Ours Theirs  
↓ Y

# Comparing sentences vs Parikh *et al.* 2016 + ELMo

Does the sentence above imply the one below, and vice versa?

*Alice carried the chair*  
*Alice held the chair*

Ground Ours Theirs  
↓ Y Y

# Comparing sentences vs Parikh *et al.* 2016 + ELMo

Does the sentence above imply the one below, and vice versa?

*Alice carried the chair*  
*Alice held the chair*

	Ground	Ours	Theirs
↓	Y	Y	N

# Comparing sentences vs Parikh *et al.* 2016 + ELMo

Does the sentence above imply the one below, and vice versa?

*Alice carried the chair*  
*Alice held the chair*

	Ground	Ours	Theirs
↓	Y	Y	N
↑	N		

# Comparing sentences vs Parikh *et al.* 2016 + ELMo

Does the sentence above imply the one below, and vice versa?

*Alice carried the chair*  
*Alice held the chair*

	Ground	Ours	Theirs
↓	Y	Y	N
↑	N	N	

# Comparing sentences vs Parikh *et al.* 2016 + ELMo

Does the sentence above imply the one below, and vice versa?

*Alice carried the chair*  
*Alice held the chair*

	Ground	Ours	Theirs
↓	Y	Y	N
↑	N	N	Y

# Comparing sentences vs Parikh *et al.* 2016 + ELMo

Does the sentence above imply the one below, and vice versa?

	Ground	Ours	Theirs
<i>Alice carried the chair</i>		Y	N
<i>Alice held the chair</i>		Y	N
<i>Alice carried the chair towards Ben</i>	↓	Y	N
<i>Alice approached Ben</i>	↑	N	Y

# Comparing sentences vs Parikh *et al.* 2016 + ELMo

Does the sentence above imply the one below, and vice versa?

	Ground	Ours	Theirs
<i>Alice carried the chair</i>			
<i>Alice held the chair</i>	↓ Y	Y	N
	↑ N	N	Y
<i>Alice carried the chair towards Ben</i>			
<i>Alice approached Ben</i>	↓ Y	Y	N
	↑ N	N	Y?



# Comparing sentences vs Parikh *et al.* 2016 + ELMo

Does the sentence above imply the one below, and vice versa?

		Ground	Ours	Theirs
<i>Alice carried the chair</i>	↓	Y	Y	N
<i>Alice held the chair</i>	↑	N	N	Y
<i>Alice carried the chair towards Ben</i>	↓	Y	Y	N
<i>Alice approached Ben</i>	↑	N	N	Y?
<i>Alice carried the chair towards Ben</i>	↓	N	N	Y?
<i>Alice left Ben</i>	↑	N	N	Y?

# Comparing sentences vs Parikh *et al.* 2016 + ELMo

Does the sentence above imply the one below, and vice versa?

		Ground	Ours	Theirs
<i>Alice carried the chair</i>	↓	Y	Y	N
<i>Alice held the chair</i>	↑	N	N	Y
<i>Alice carried the chair towards Ben</i>	↓	Y	Y	N
<i>Alice approached Ben</i>	↑	N	N	Y?
<i>Alice carried the chair towards Ben</i>	↓	N	N	Y?
<i>Alice left Ben</i>	↑	N	N	Y?
<i>Alice picked up the chair, and Ben put down the bag</i>	↓	N	N	Y
<i>Ben picked up the chair, and Alice put down the bag</i>	↑	N	N	Y

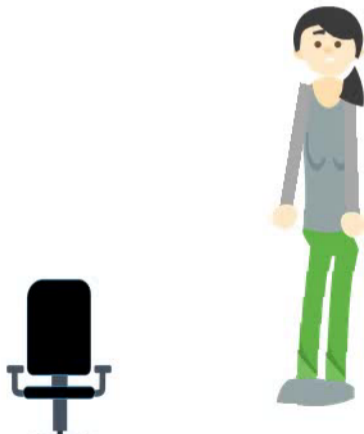
# What about other languages?

## What about other languages?

爱丽丝靠近了一把椅子。

## What about other languages?

爱丽丝靠近了一把椅子。



# Word Embedding Association Test (WEAT)

# Word Embedding Association Test (WEAT)

*Women* = {Anna, Mary}

*Men* = {Dave, John}

*Work* = {office, desk}

*Home* = {children, home}

# Word Embedding Association Test (WEAT)

*Women* = {Anna, Mary}

*Men* = {Dave, John}

Targets

*Work* = {office, desk}

*Home* = {children, home}

Attributes



# Word Embedding Association Test (WEAT)

*Women* = {Anna, Mary}

*Men* = {Dave, John}

Targets

*Work* = {office, desk}

*Home* = {children, home}

Attributes

Compare the distances between *Women* *Work* and *Home*  
and between *Men* *Work* and *Home*

# Word Embedding Association Test (WEAT)

*Women* = {Anna, Mary}

*Men* = {Dave, John}

Targets

*Work* = {office, desk}

*Home* = {children, home}

Attributes

Compare the distances between *Women* *Work* and *Home*  
and between *Men* *Work* and *Home*

If they are very different then there is some bias  
because this has practical consequences for the inferences that networks make.  
There are variants like SEAT which test whole sentences.

# Grounded WEAT

Dave



John



Steve



Anna



Mary



Beth



doctor



lawyer



police



librarian



teacher



secretary



doctor



lawyer



police



librarian



teacher



secretary



# Different groupings answer different questions

# Different groupings answer different questions

3 tests: word, sentence, and context  
6 gender bias tests and 7 racial bias tests  
4 popular models (VisualBERT, VL-BERT, ViLBERT, LXMERT)

# Different groupings answer different questions

3 tests: word, sentence, and context  
6 gender bias tests and 7 racial bias tests  
4 popular models (VisualBERT, VL-BERT, ViLBERT, LXMERT)

Do multimodal models have social biases?

# Different groupings answer different questions

3 tests: word, sentence, and context

6 gender bias tests and 7 racial bias tests

4 popular models (VisualBERT, VL-BERT, ViLBERT, LXMERT)

Do multimodal models have social biases? They all do

# Different groupings answer different questions

3 tests: word, sentence, and context  
6 gender bias tests and 7 racial bias tests  
4 popular models (VisualBERT, VL-BERT, ViLBERT, LXMERT)

Do multimodal models have social biases? They all do  
Can counterstereotypical visual evidence offset a bias?



# Different groupings answer different questions

3 tests: word, sentence, and context

6 gender bias tests and 7 racial bias tests

4 popular models (VisualBERT, VL-BERT, ViLBERT, LXMERT)

Do multimodal models have social biases? They all do

Can counterstereotypical visual evidence offset a bias? No

Do biases come from vision or language?

# Different groupings answer different questions

3 tests: word, sentence, and context

6 gender bias tests and 7 racial bias tests

4 popular models (VisualBERT, VL-BERT, ViLBERT, LXMERT)

Do multimodal models have social biases? They all do

Can counterstereotypical visual evidence offset a bias? No

Do biases come from vision or language? ViLBERT clearly language, mix for the others

# A missing algorithm in our language learning story

# A missing algorithm in our language learning story



## A missing algorithm in our language learning story



Overwhelm with realistic but synthetic visual evidence against biases.

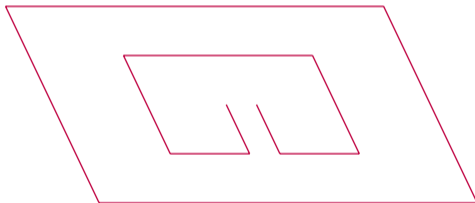
# Communication isn't all verbal



# Communication isn't all verbal

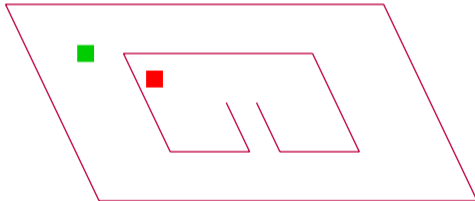


## Sampling-based planning with language

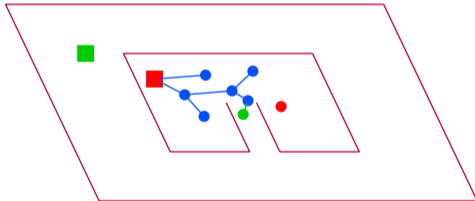




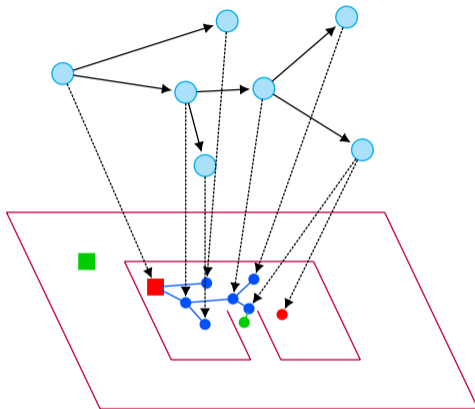
## Sampling-based planning with language



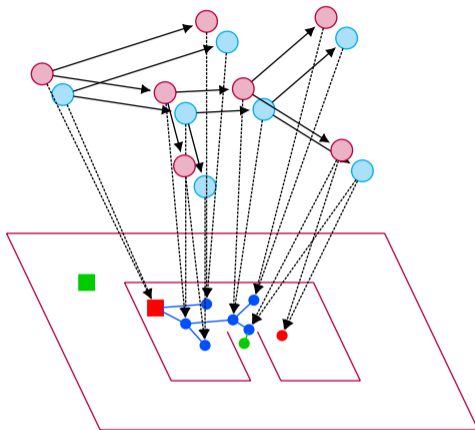
# Sampling-based planning with language



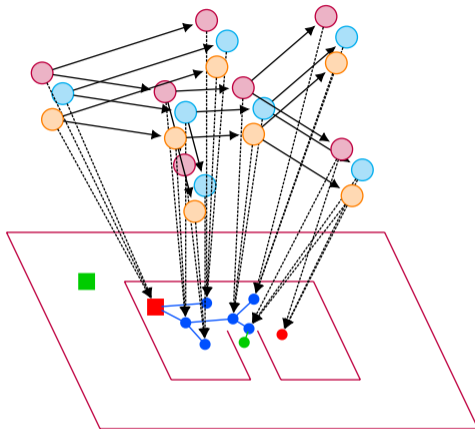
# Sampling-based planning with language



# Sampling-based planning with language

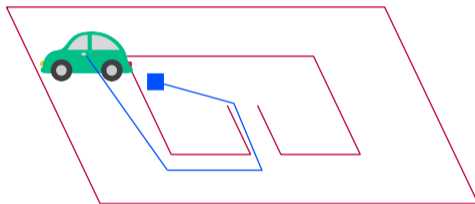


# Sampling-based planning with language



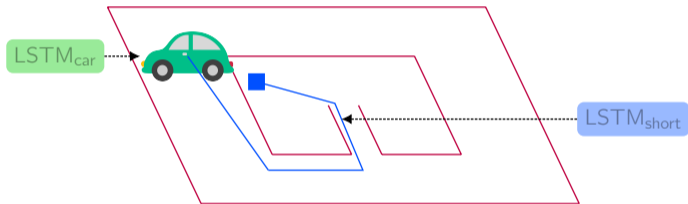


# Planning and language



Go to the **car** quickly.

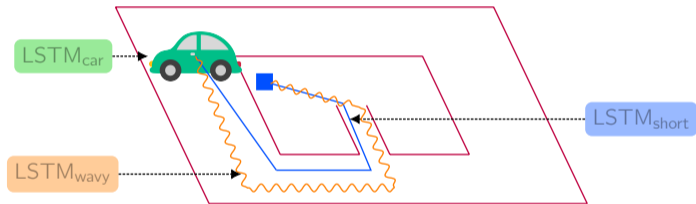
# Planning and language



Go to the **car** quickly.

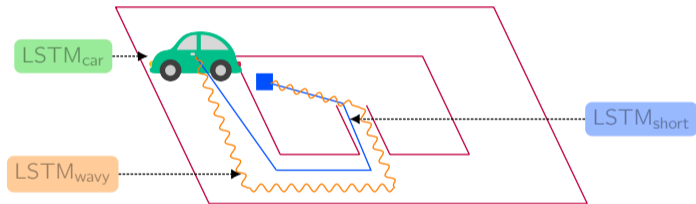


# Planning and language



Weave to the green car.

# Planning and language



Weave to the green car.

# Turn parses into networks that encode sentences

# Turn parses into networks that encode sentences

Pick up the black triangle below the orange ball.

# Turn parses into networks that encode sentences

Pick up the black triangle below the orange ball.



# Turn parses into networks that encode sentences

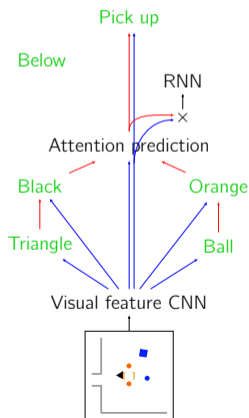
Pick up the black triangle below the orange ball.

Visual feature CNN



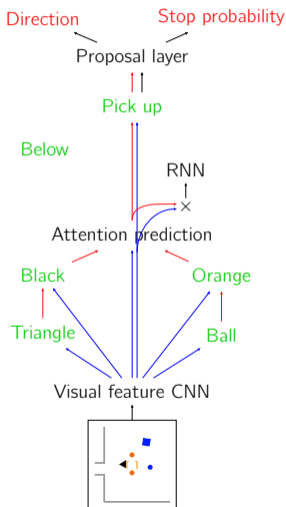
# Turn parses into networks that encode sentences

Pick up the black triangle below the orange ball.



# Turn parses into networks that encode sentences

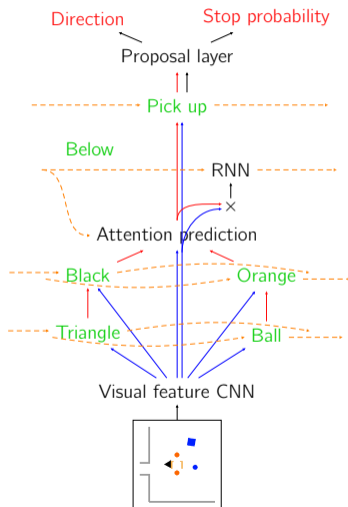
Pick up the black triangle below the orange ball.





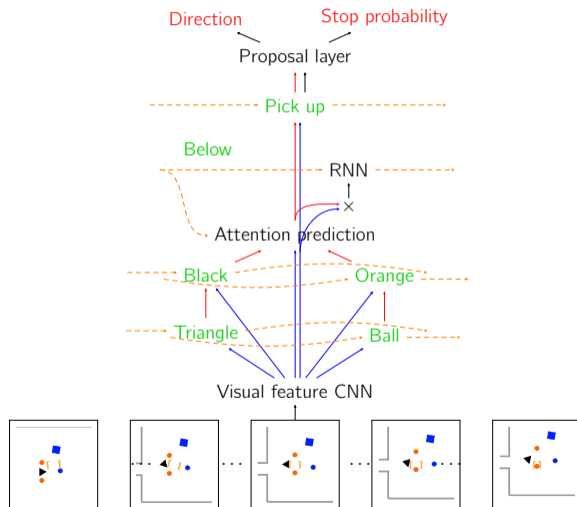
# Turn parses into networks that encode sentences

Pick up the black triangle below the orange ball.



# Turn parses into networks that encode sentences

Pick up the black triangle below the orange ball.



# Compositionality matters, but how is unclear

Planner accuracy for different models

Same number of parameters, same implementation, same optimizer, same hyperparameters

# Compositionality matters, but how is unclear

## Planner accuracy for different models

Same number of parameters, same implementation, same optimizer, same hyperparameters

Dataset	RNN	Compositional RNN	Compositional RNN (bad tree)
gSCAN in domain	97%	96%	96%
gSCAN out of domain	54%	96%	58%
gSCAN target length 15	95%	93%	
gSCAN target length 16	19%	91%	
gSCAN target length 17	1%	88%	
gSCAN target length 18	≈0%	57%	

# Compositionality matters, but how is unclear

## Planner accuracy for different models

Same number of parameters, same implementation, same optimizer, same hyperparameters

Dataset	RNN	Compositional RNN	Compositional RNN (bad tree)
gSCAN in domain	97%	96%	96%
gSCAN out of domain	54%	96%	58%
gSCAN target length 15	95%	93%	
gSCAN target length 16	19%	91%	
gSCAN target length 17	1%	88%	
gSCAN target length 18	$\approx 0\%$	57%	

Our compositional network is robust to new combinations

But why are we robust to longer sequences?

# What theories can we generalize to?

Pick up the box.

Propositional Logic

# What theories can we generalize to?

Pick up the box.  
Pick up all the boxes.

Propositional Logic  
FOL

# What theories can we generalize to?

Pick up the box.

Pick up all the boxes.

Every time a box falls on the ground, pick it up.

Propositional Logic

FOL

LTL, STL, etc.



# What theories can we generalize to?

Pick up the box.

Pick up all the boxes.

Every time a box falls on the ground, pick it up.

If the box is stacked precariously, fix it.

Propositional Logic

FOL

LTL, STL, etc.

Physics

# What theories can we generalize to?

Pick up the box.

Pick up all the boxes.

Every time a box falls on the ground, pick it up.

If the box is stacked precariously, fix it.

If someone wants to drop the box, stop them.

Propositional Logic

FOL

LTL, STL, etc.

Physics

Inverse planning, ToM

# What theories can we generalize to?

Pick up the box.

Pick up all the boxes.

Every time a box falls on the ground, pick it up.

If the box is stacked precariously, fix it.

If someone wants to drop the box, stop them.

If the box is about to fall, catch it.

Propositional Logic

FOL

LTL, STL, etc.

Physics

Inverse planning, ToM

Possibility, Modal logic

# What theories can we generalize to?

Pick up the box.

Pick up all the boxes.

Every time a box falls on the ground, pick it up.

If the box is stacked precariously, fix it.

If someone wants to drop the box, stop them.

If the box is about to fall, catch it.

Pick up the biggest box.

Propositional Logic

FOL

LTL, STL, etc.

Physics

Inverse planning, ToM

Possibility, Modal logic

Scalar implicatures

# What theories can we generalize to?

Pick up the box.

Pick up all the boxes.

Every time a box falls on the ground, pick it up.

If the box is stacked precariously, fix it.

If someone wants to drop the box, stop them.

If the box is about to fall, catch it.

Pick up the biggest box.

Get the box that won't leak.

Propositional Logic

FOL

LTL, STL, etc.

Physics

Inverse planning, ToM

Possibility, Modal logic

Scalar implicatures

Modification

# What theories can we generalize to?

Pick up the box.

Pick up all the boxes.

Every time a box falls on the ground, pick it up.

If the box is stacked precariously, fix it.

If someone wants to drop the box, stop them.

If the box is about to fall, catch it.

Pick up the biggest box.

Get the box that won't leak.

Show me the shiny side.

Propositional Logic

FOL

LTL, STL, etc.

Physics

Inverse planning, ToM

Possibility, Modal logic

Scalar implicatures

Modification

Grounding

# What theories can we generalize to?

Pick up the box.

Pick up all the boxes.

Every time a box falls on the ground, pick it up.

If the box is stacked precariously, fix it.

If someone wants to drop the box, stop them.

If the box is about to fall, catch it.

Pick up the biggest box.

Get the box that won't leak.

Show me the shiny side.

Be friendly

Propositional Logic

FOL

LTL, STL, etc.

Physics

Inverse planning, ToM

Possibility, Modal logic

Scalar implicatures

Modification

Grounding

Social interactions

# What theories can we generalize to?

Pick up the box.

Pick up all the boxes.

Every time a box falls on the ground, pick it up.

If the box is stacked precariously, fix it.

If someone wants to drop the box, stop them.

If the box is about to fall, catch it.

Pick up the biggest box.

Get the box that won't leak.

Show me the shiny side.

Be friendly

...

Propositional Logic

FOL

LTL, STL, etc.

Physics

Inverse planning, ToM

Possibility, Modal logic

Scalar implicatures

Modification

Grounding

Social interactions



# What theories can we generalize to?

Pick up the box.

Pick up all the boxes.

Every time a box falls on the ground, pick it up.

If the box is stacked precariously, fix it.

If someone wants to drop the box, stop them.

If the box is about to fall, catch it.

Pick up the biggest box.

Get the box that won't leak.

Show me the shiny side.

Be friendly

Propositional Logic

FOL

LTL, STL, etc.

Physics

Inverse planning, ToM

Possibility, Modal logic

Scalar implicatures

Modification

Grounding

Social interactions

...

You can't just hope to generalize between these domains!

But in the real world . . .



But in the real world . . .



# Vision is integral to grounded language acquisition

Torralba & Efros, 2011; Berzak *et al.*, 2017

# Vision is integral to grounded language acquisition

Machine performance on ImageNet is around 97%

# Vision is integral to grounded language acquisition

Machine performance on ImageNet is around 97%  
Human-level performance on ImageNet is around 94%

# Vision is integral to grounded language acquisition

Machine performance on ImageNet is around 97%  
Human-level performance on ImageNet is around 94%  
Machines outperform humans according standard metrics!

# Vision is integral to grounded language acquisition

Machine performance on ImageNet is around 97%

Human-level performance on ImageNet is around 94%

Machines outperform humans according standard metrics!

Performance on datasets is not predictive of real-world performance.



# Vision is integral to grounded language acquisition

Machine performance on ImageNet is around 97%

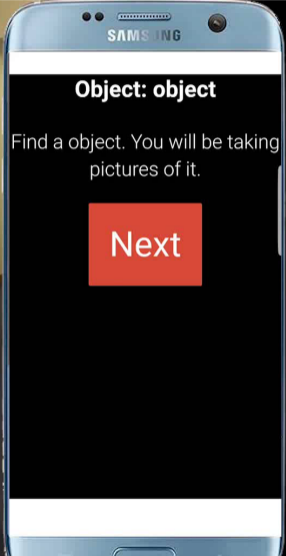
Human-level performance on ImageNet is around 94%

Machines outperform humans according standard metrics!

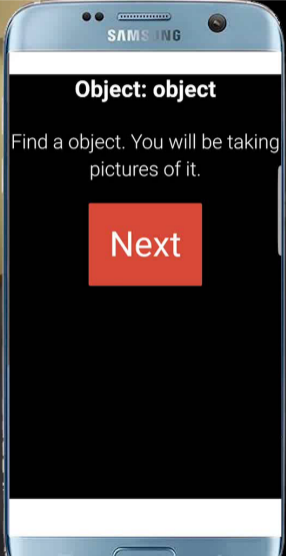
Performance on datasets is not predictive of real-world performance.

Or even performance on other datasets!

# Data collection in action



# Data collection in action



313 object classes

313 object classes

50k images

313 object classes

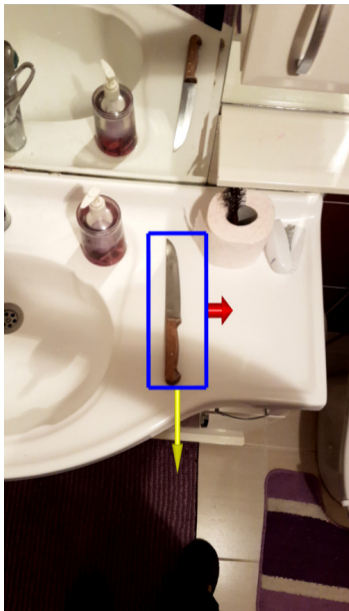
50k images

No training set!

313 object classes

50k images

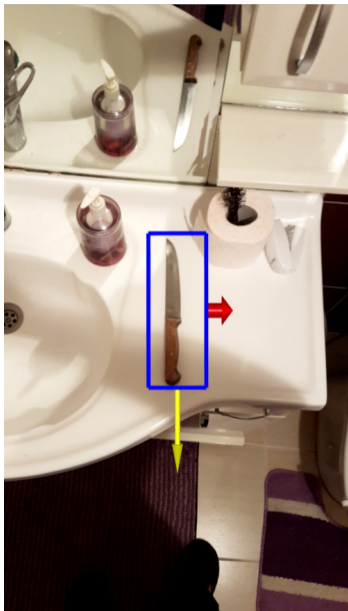
No training set!



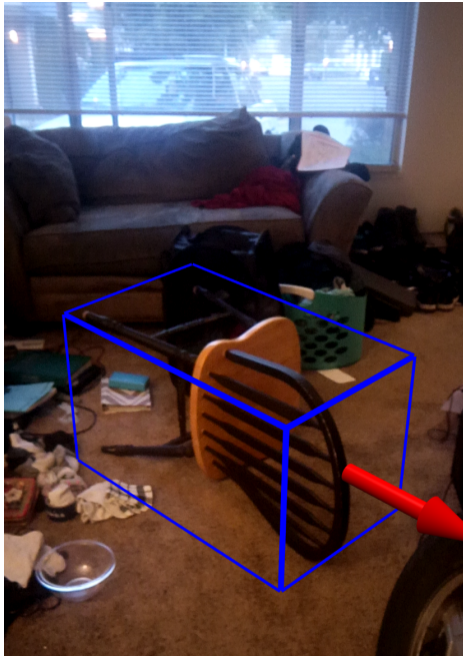
313 object classes

50k images

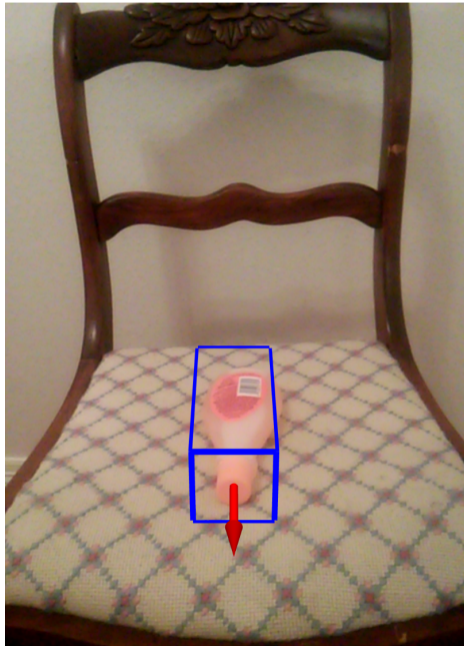
No training set!

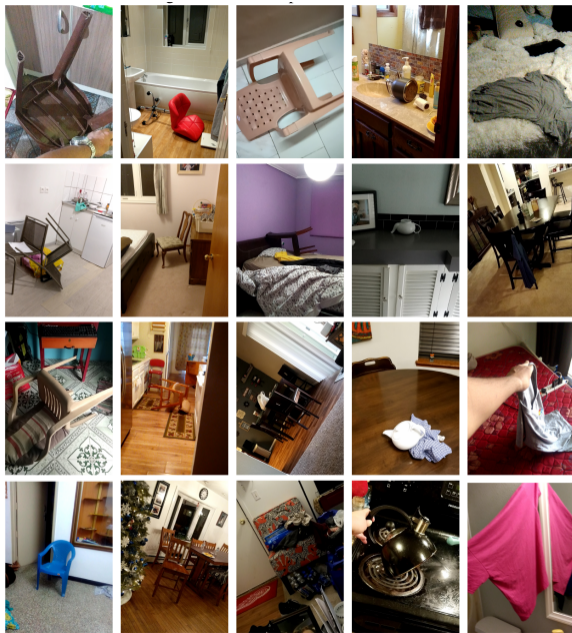


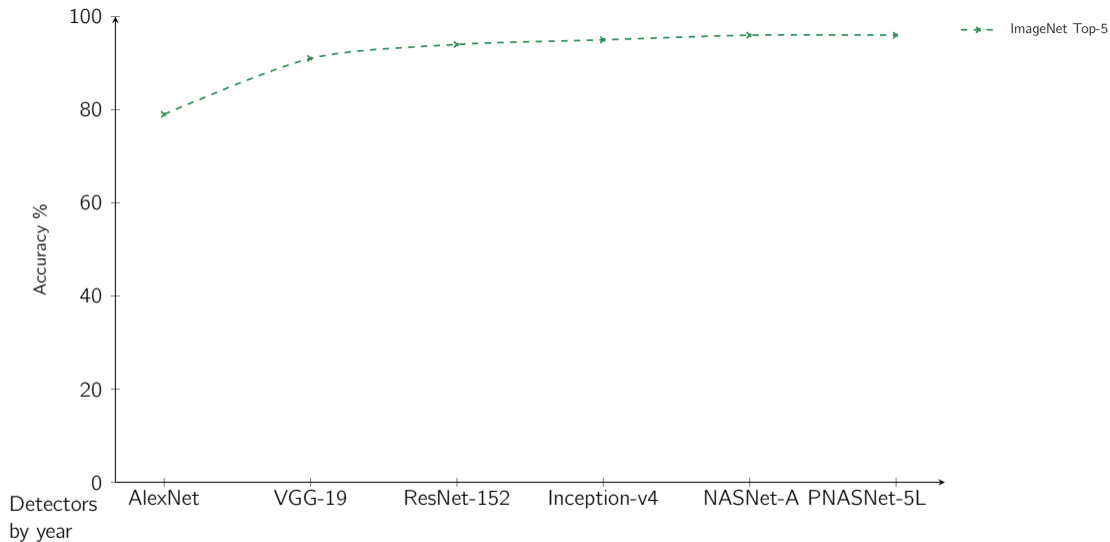


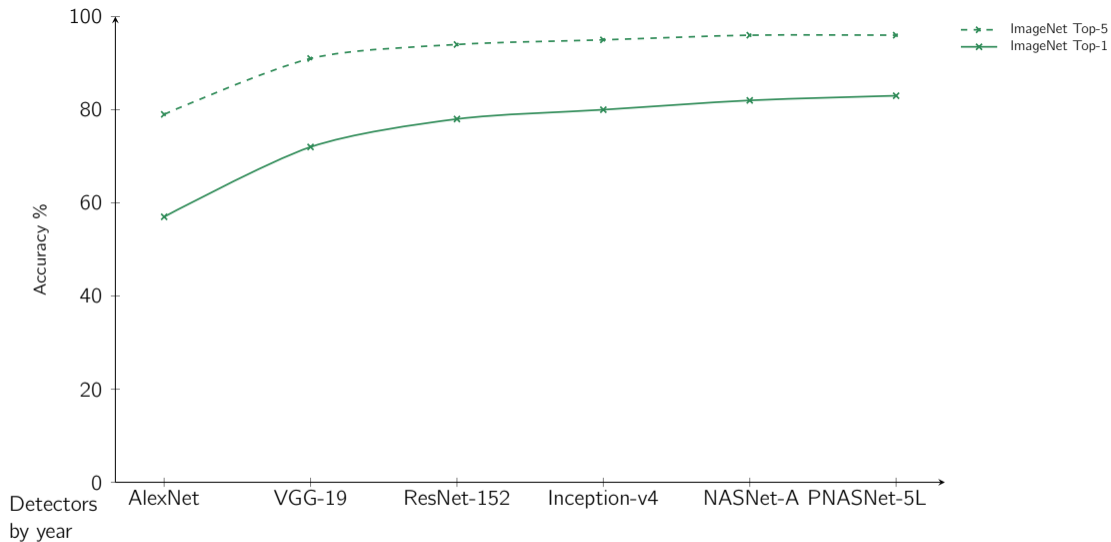


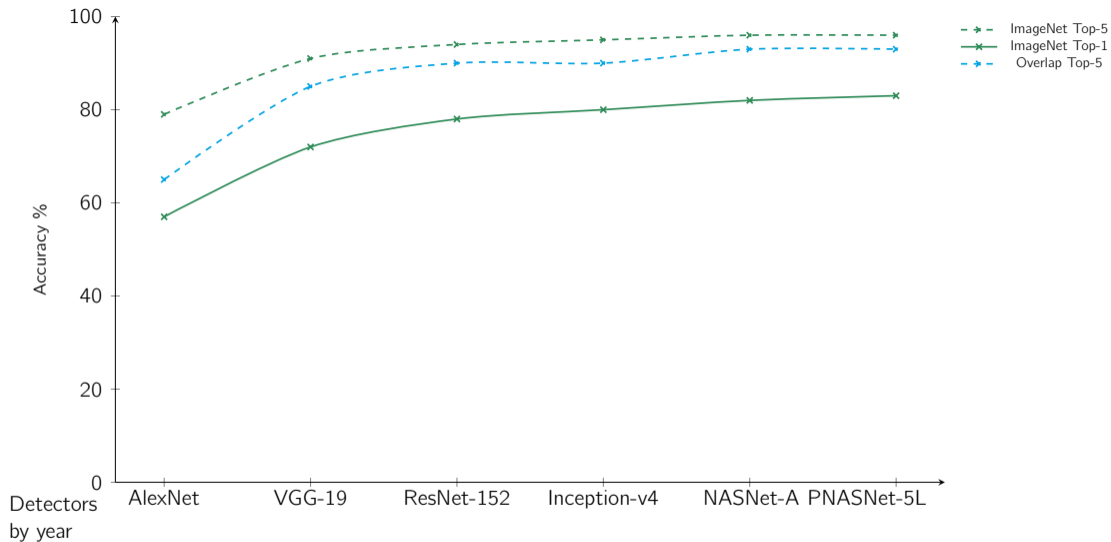


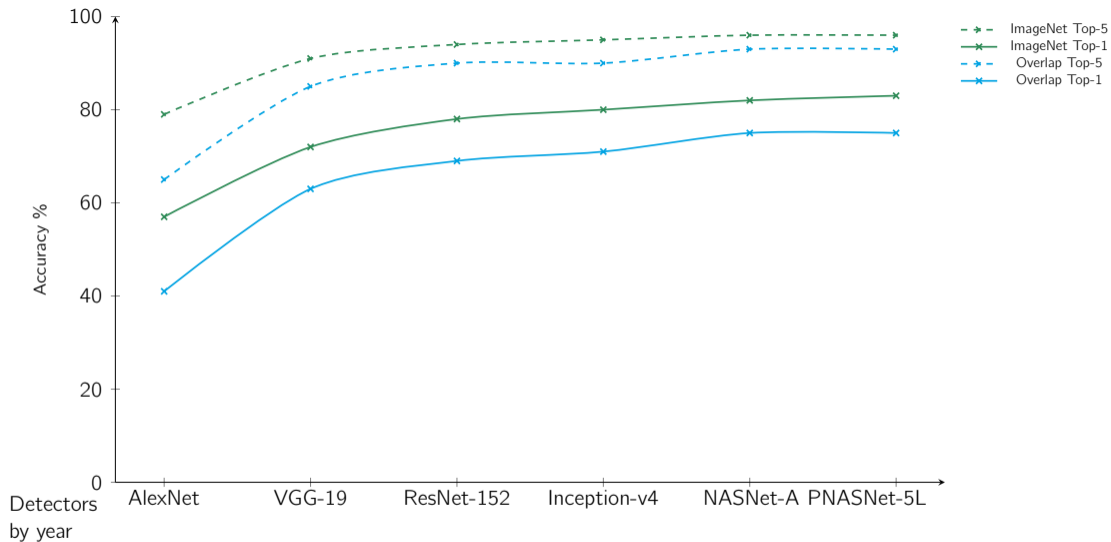




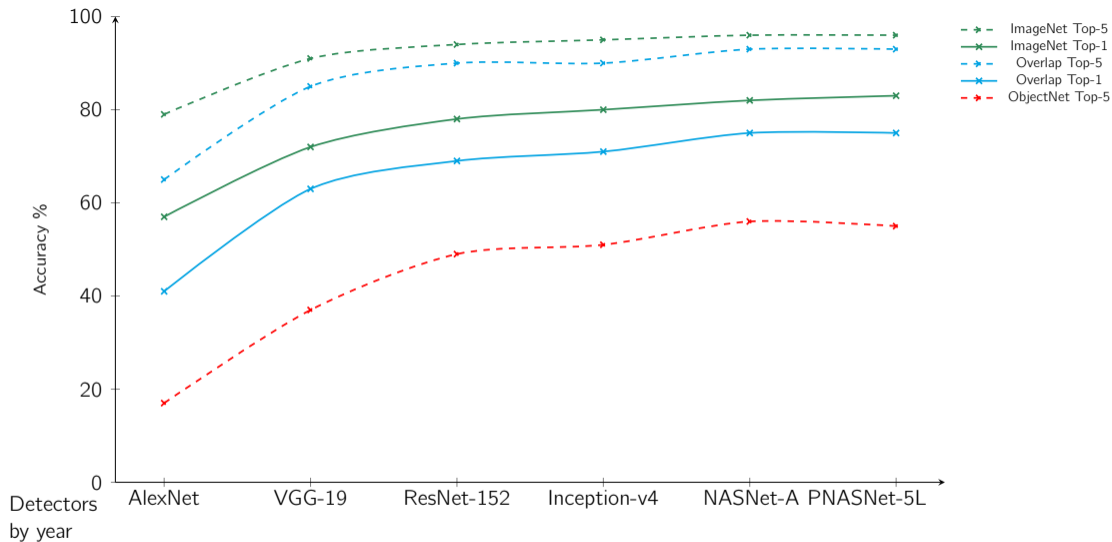


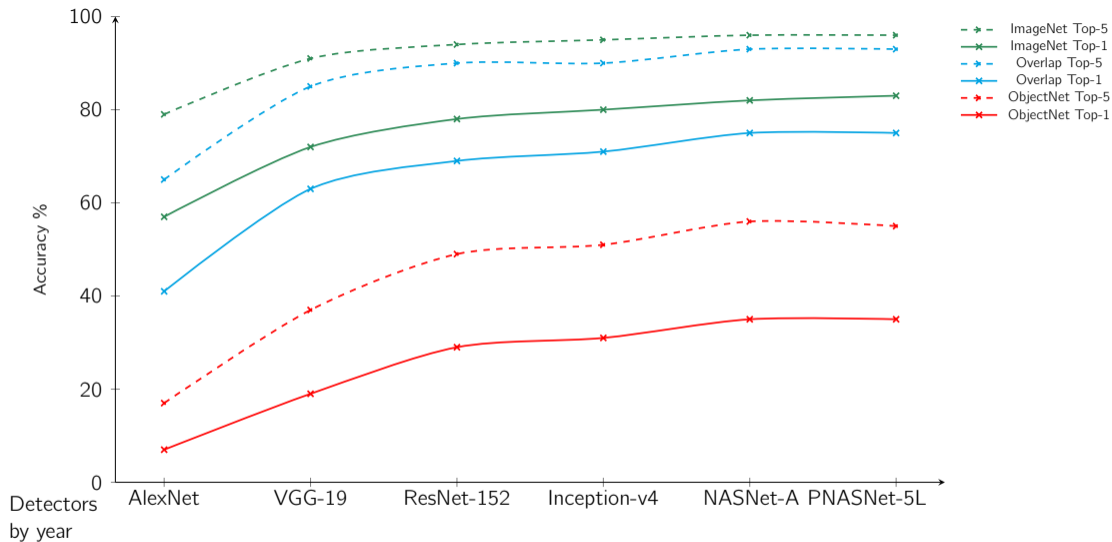


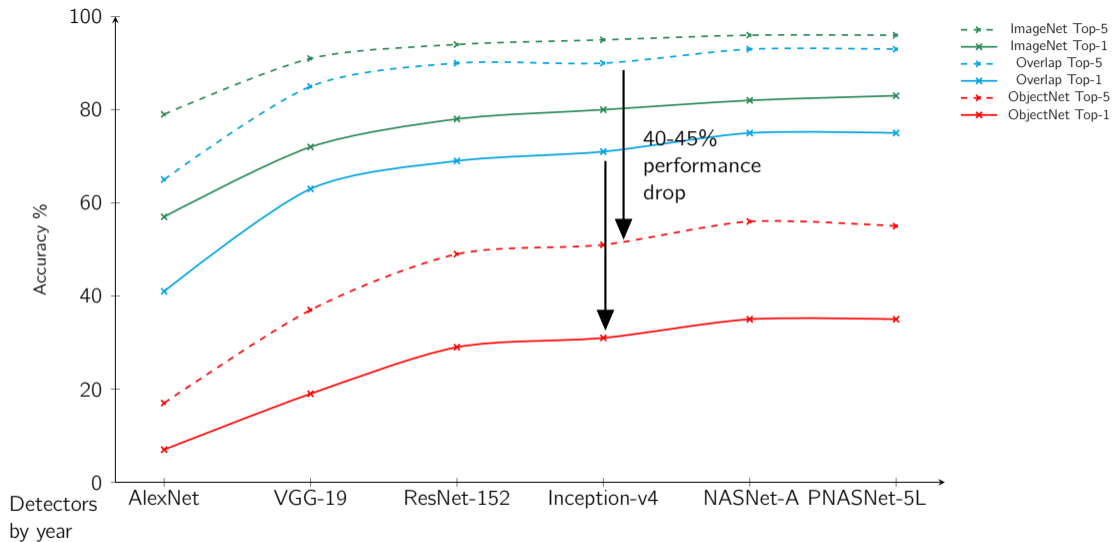






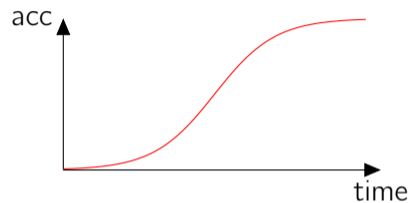




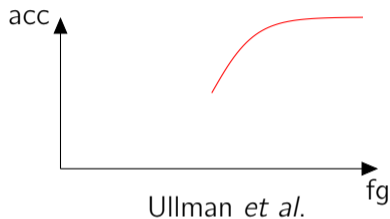
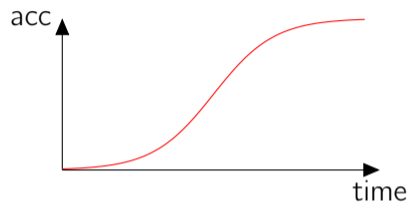


# Next steps in understanding human vision

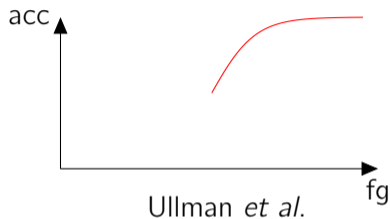
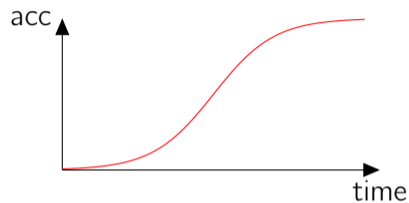
## Next steps in understanding human vision



## Next steps in understanding human vision

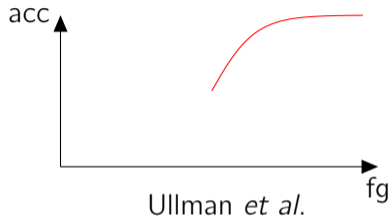
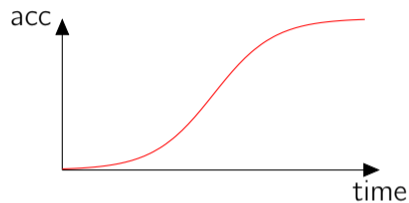


## Next steps in understanding human vision



Accuracy(Foreground, Background, Time)

## Next steps in understanding human vision

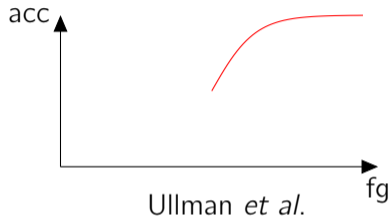
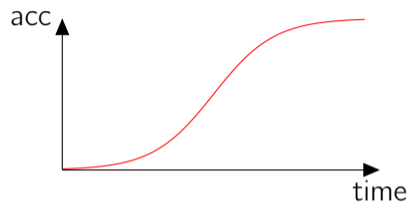


Accuracy(Foreground, Background, Time)

What object features (shape, texture, etc.) lead to similar response curves?



## Next steps in understanding human vision

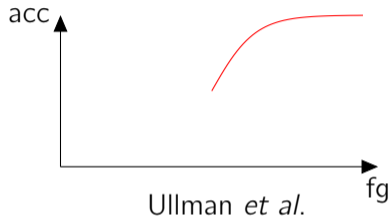
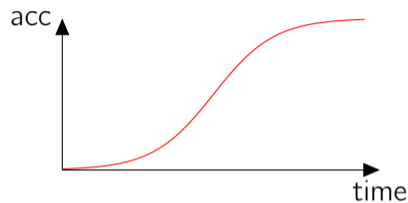


Accuracy(Foreground, Background, Time)

What object features (shape, texture, etc.) lead to similar response curves?

What are the parameters of these curves?

## Next steps in understanding human vision



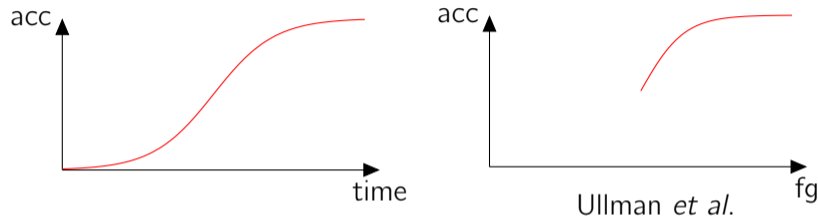
Accuracy(Foreground, Background, Time)

What object features (shape, texture, etc.) lead to similar response curves?

What are the parameters of these curves?

Are there any other discontinuities?

## Next steps in understanding human vision



Accuracy(Foreground, Background, Time)

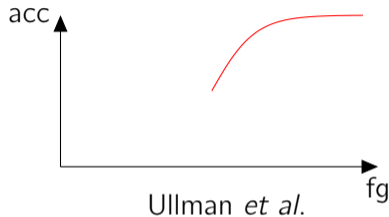
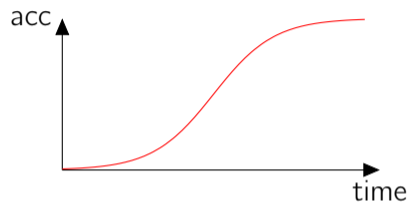
What object features (shape, texture, etc.) lead to similar response curves?

What are the parameters of these curves?

Are there any other discontinuities?

Can we find mode switches? (feedforward vs. feedback)

## Next steps in understanding human vision



Ullman *et al.*

Accuracy(Foreground, Background, Time)

What object features (shape, texture, etc.) lead to similar response curves?

What are the parameters of these curves?

Are there any other discontinuities?

Can we find mode switches? (feedforward vs. feedback)

Can human data constrain networks?

## Captioning datasets are really biased . . .

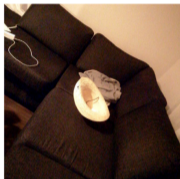
Spoken ObjectNet to address this and to test out of domain generalization.

# Captioning datasets are really biased . . .

Spoken ObjectNet to address this and to test out of domain generalization.



a jar of honey being held on its side of the yellow lid in front of old flowers in a vase on a bathroom sink countertop



a straw like sun hat upside down on a brown sectional with a gray piece of clothing in the corner



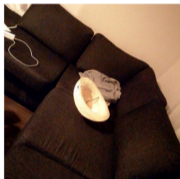
someone holding a bottle of extra virgin olive oil upside down over a white washer with a plug-in behind it

# Captioning datasets are really biased . . .

Spoken ObjectNet to address this and to test out of domain generalization.



a jar of honey being held on its side of the yellow lid in front of old flowers in a vase on a bathroom sink countertop



a straw like sun hat upside down on a brown sectional with a gray piece of clothing in the corner



someone holding a bottle of extra virgin olive oil upside down over a white washer with a plug-in behind it

On a retrieval task trained with Places-400k R@10 goes from 0.735 to 0.118!

A larger drop than for object recognition on ObjectNet

Palmer, Rouditchenko, Barbu, Katz, Glass in review

# Big data neuroscience for language

Yaari, Singh, Cases, Subramaniam, Katz, Kreiman, Barbu in prep.



# Big data neuroscience for language

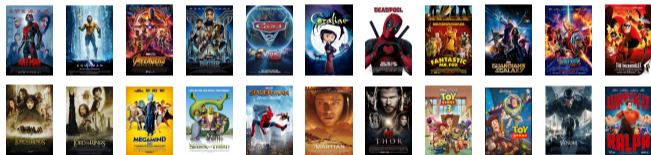
Yaari, Singh, Cases, Subramaniam, Katz, Kreiman, Barbu in prep.

# Big data neuroscience for language



Yaari, Singh, Cases, Subramaniam, Katz, Kreiman, Barbu in prep.

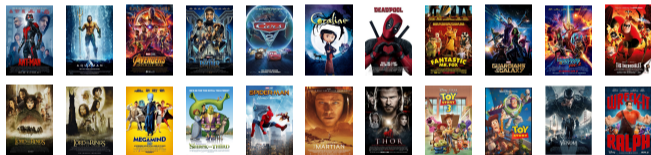
# Big data neuroscience for language



Collected  $\approx 40$  hours of data while 7 subjects watched movies.

Yaari, Singh, Cases, Subramaniam, Katz, Kreiman, Barbu in prep.

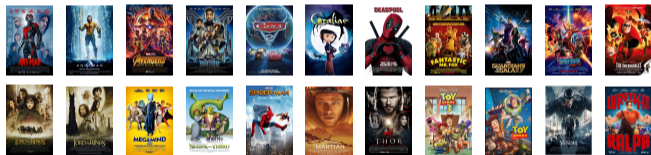
# Big data neuroscience for language



Collected  $\approx 40$  hours of data while 7 subjects watched movies.  
 $153 \pm 26.6$  electrodes per subject, 5 male / 2 female, 12.5 years old  $\pm 5.23$

Yaari, Singh, Cases, Subramaniam, Katz, Kreiman, Barbu in prep.

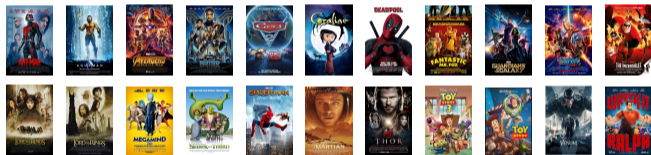
# Big data neuroscience for language



Collected  $\approx 40$  hours of data while 7 subjects watched movies.  
 $153 \pm 26.6$  electrodes per subject, 5 male / 2 female, 12.5 years old  $\pm 5.23$   
1,079 electrodes, 27,981 sentences, 169,314 words

Yaari, Singh, Cases, Subramaniam, Katz, Kreiman, Barbu in prep.

# Big data neuroscience for language



Collected  $\approx 40$  hours of data while 7 subjects watched movies.  
153  $\pm$  26.6 electrodes per subject, 5 male / 2 female, 12.5 years old  $\pm$  5.23  
1,079 electrodes, 27,981 sentences, 169,314 words  
10-100x more data per subject than previous language datasets.

Yaari, Singh, Cases, Subramaniam, Katz, Kreiman, Barbu in prep.

## Some levels of linguistic analysis

*I don't think Gandalf meant for us to come this way.*

# Some levels of linguistic analysis

Word segmentation

| I | don't | think | Gandalf | meant | for | us | to | come | this | way. |



# Some levels of linguistic analysis

Word segmentation

| I | don't | think | Gandalf | meant | for | us | to | come | this | way. |

Phonemes

aɪ doʊnt θɪŋk 'gændɔlf mənt fər əs tə kʌm ðɪs weɪ

# Some levels of linguistic analysis

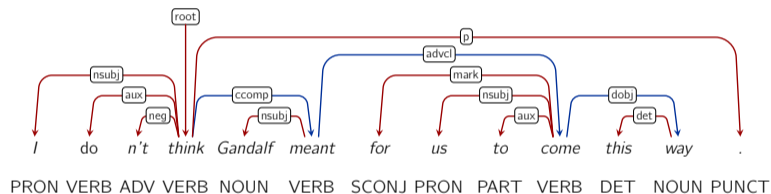
Word segmentation

| I | don't | think | Gandalf | meant | for | us | to | come | this | way | . |

Phonemes

aɪ dɒnt θɪŋk ˈgændɔlf mənt fər əs tə kʌm ðɪs weɪ

Parse structure



# Some levels of linguistic analysis

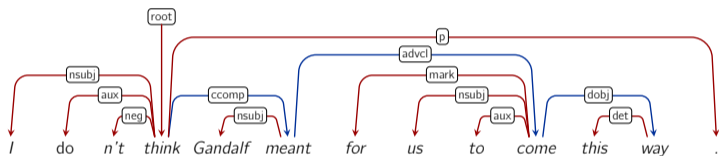
Word segmentation

| I | don't | think | Gandalf | meant | for | us | to | come | this | way | .

Phonemes

aɪ dɒnt θɪŋk ˈɡændɔlf mənt fər əs tə kʌm ðɪs weɪ

Parse structure



Part of speech

PRON VERB ADV VERB NOUN VERB SCONJ PRON PART VERB DET NOUN PUNCT

# Some levels of linguistic analysis

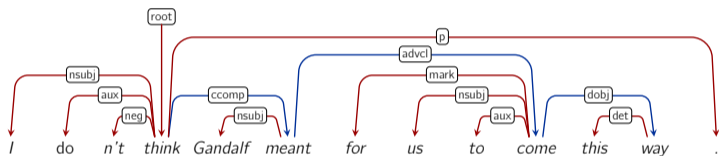
Word segmentation

| I | don't | think | Gandalf | meant | for | us | to | come | this | way | .

Phonemes

aɪ dɒnt θɪŋk ˈgændəlf mənt fər əs tə kʌm ðɪs weɪ

Parse structure



Part of speech

PRON VERB ADV VERB NOUN VERB SCONJ PRON PART VERB DET NOUN PUNCT



Thematic relations

# Some levels of linguistic analysis

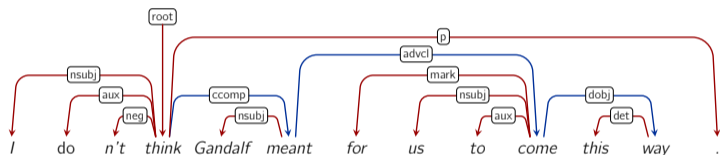
Word segmentation

| I | don't | think | Gandalf | meant | for | us | to | come | this | way | . |

Phonemes

aɪ doʊnt θɪŋk ˈɡændəlf mənt fər əs tə kʌm ðɪs weɪ

Parse structure



Part of speech

PRON VERB ADV VERB NOUN VERB SCONJ PRON PART VERB DET NOUN PUNCT



Thematic relations

Sentiment analysis

Negative

# Some levels of linguistic analysis

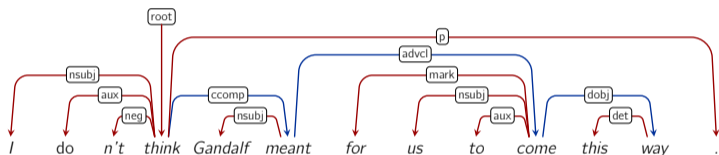
Word segmentation

| I | don't | think | Gandalf | meant | for | us | to | come | this | way | . |

Phonemes

aɪ dɒnt θɪŋk ˈgændəlf mənt fər əs tə kʌm ðɪs weɪ

Parse structure



Part of speech

PRON VERB ADV VERB NOUN VERB SCONJ PRON PART VERB DET NOUN PUNCT



Thematic relations

Sentiment analysis

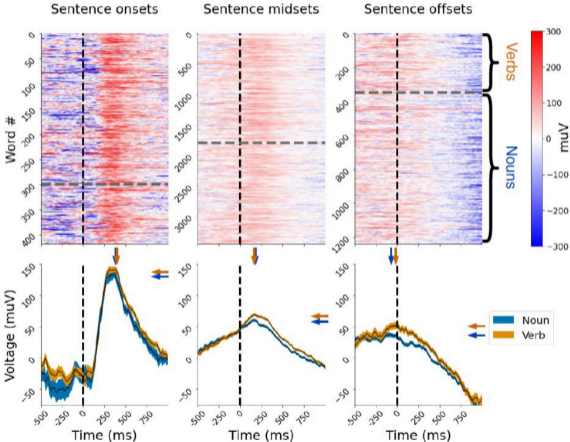
Negative

Semantics

...

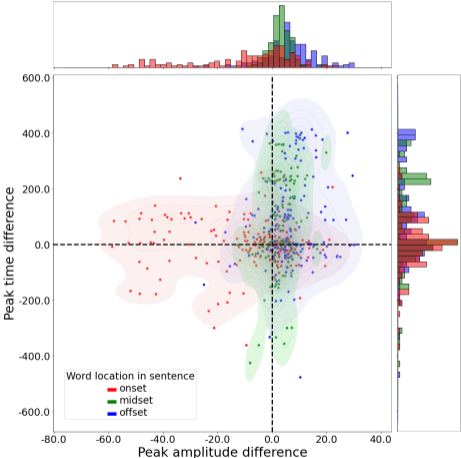
# Nouns and verbs, one electrode

Every **day** I go outside and **look** at the vast **horizons**.



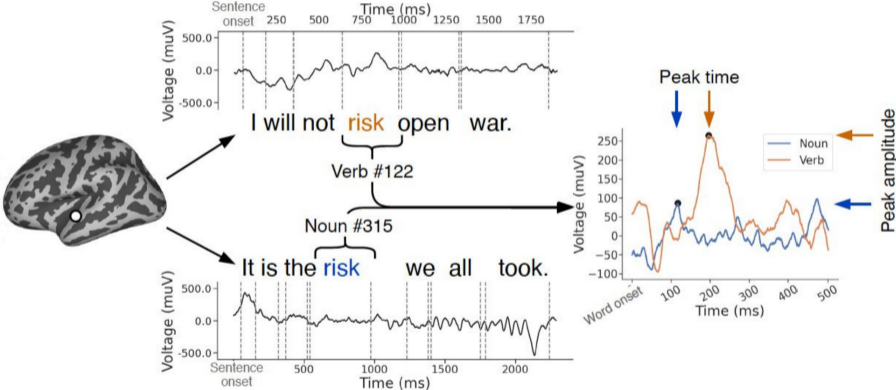
Yaari, Singh, Cases, Subramaniam, Katz, Kreiman, Barbu in prep.

# Nouns and verbs, all electrodes





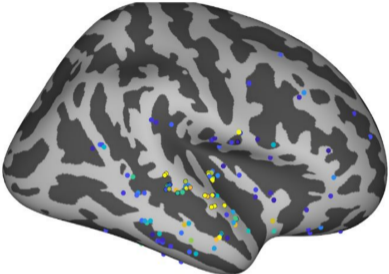
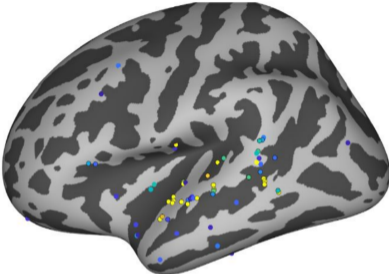
# Homonyms allow for controlled experiments



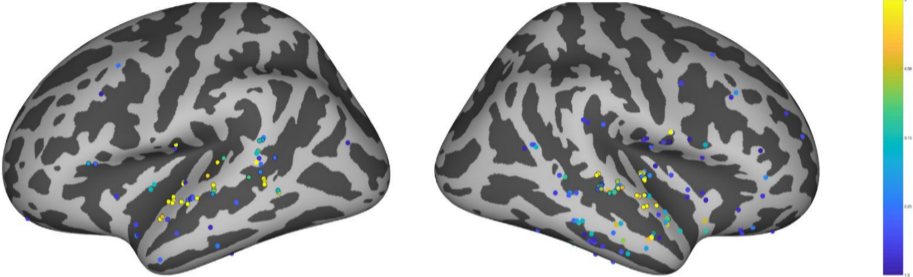
Yaari, Singh, Cases, Subramaniam, Katz, Kreiman, Barbu in prep.

# Decoding nouns and verbs — homonyms

# Decoding nouns and verbs — homonyms

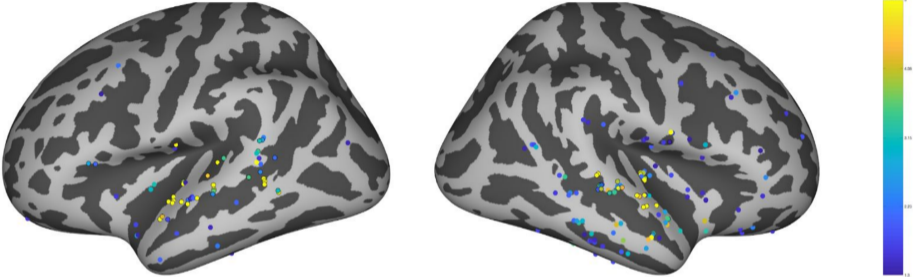


# Decoding nouns and verbs — homonyms



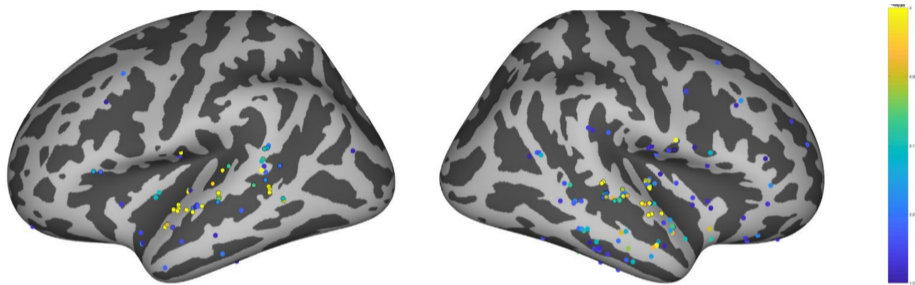
Train on nouns and verbs, hold out all homonyms

# Decoding nouns and verbs — homonyms



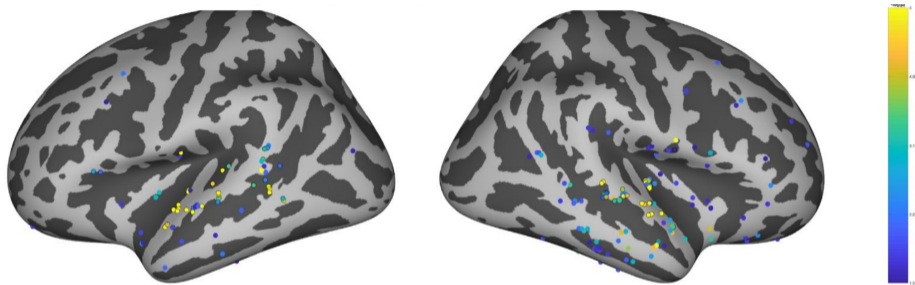
Train on nouns and verbs, hold out all homonyms  
At training time no word appears as both as noun and a verb

## Decoding nouns and verbs — homonyms



Train on nouns and verbs, hold out all homonyms  
At training time no word appears as both as noun and a verb  
At test time, test on only homonyms:

## Decoding nouns and verbs — homonyms

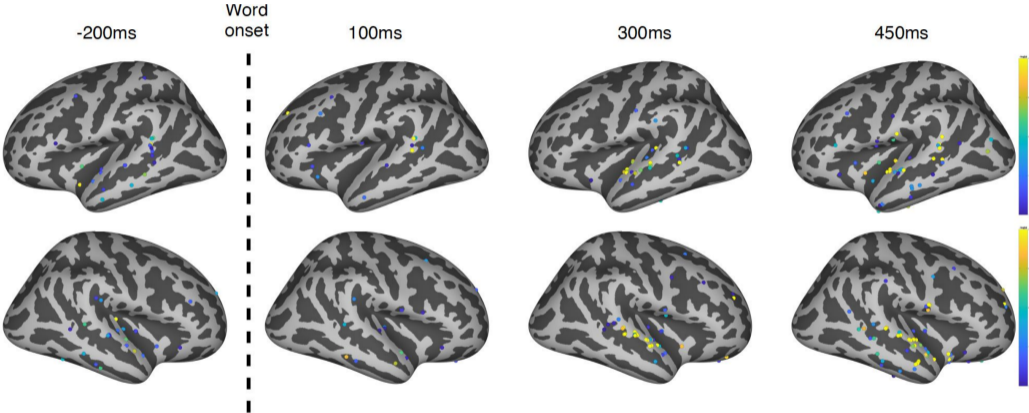


Train on nouns and verbs, hold out all homonyms  
At training time no word appears as both as noun and a verb  
At test time, test on only homonyms:  
Generalize to unseen words and unseen word-POS pairs

# Decoding nouns and verbs in time

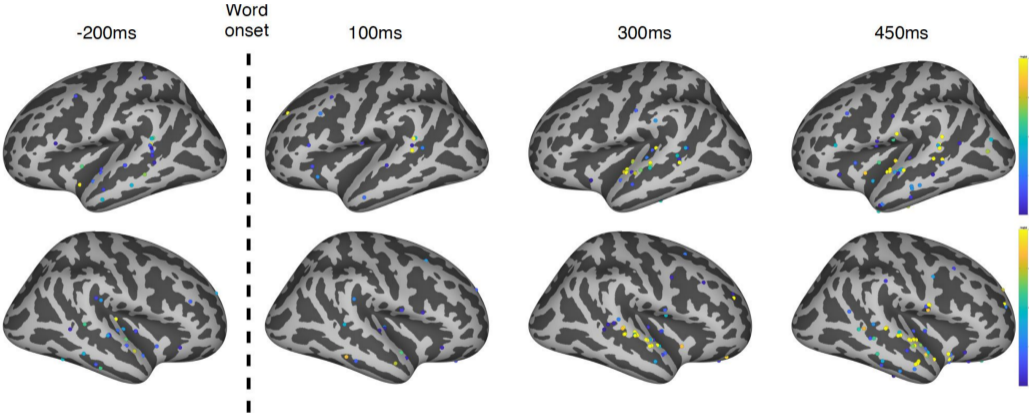


# Decoding nouns and verbs in time



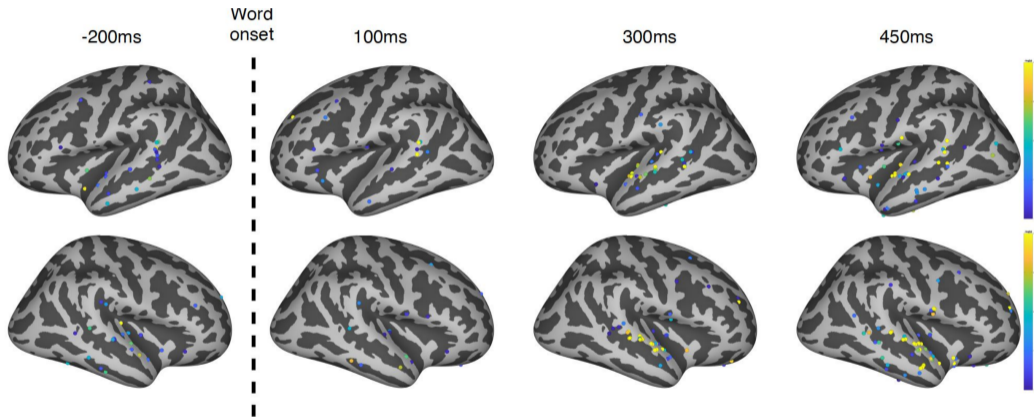
Yaari, Singh, Cases, Subramaniam, Katz, Kreiman, Barbu in prep.

# Decoding nouns and verbs in time



Well before a word is uttered you predict the POS of that word

# Decoding nouns and verbs in time



Well before a word is uttered you predict the POS of that word  
When the word is uttered predictions are updated in STG

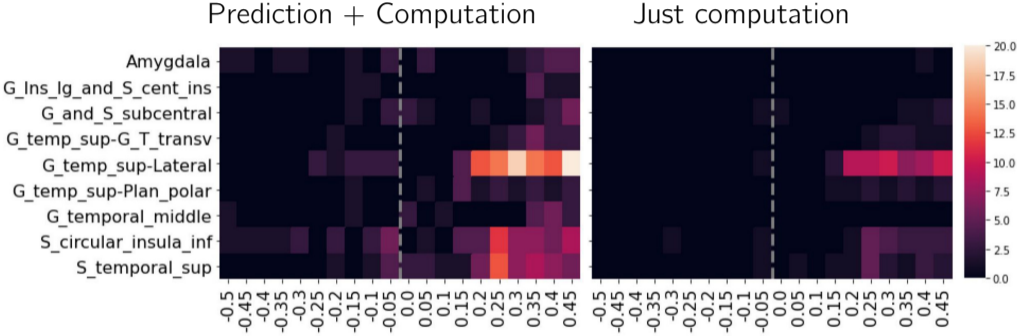
Yaari, Singh, Cases, Subramaniam, Katz, Kreiman, Barbu in prep.

# Prediction vs computation

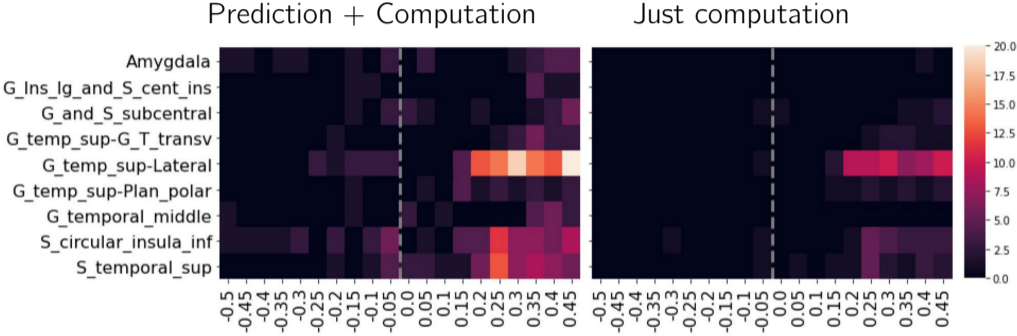
Prediction + Computation

Just computation

# Prediction vs computation

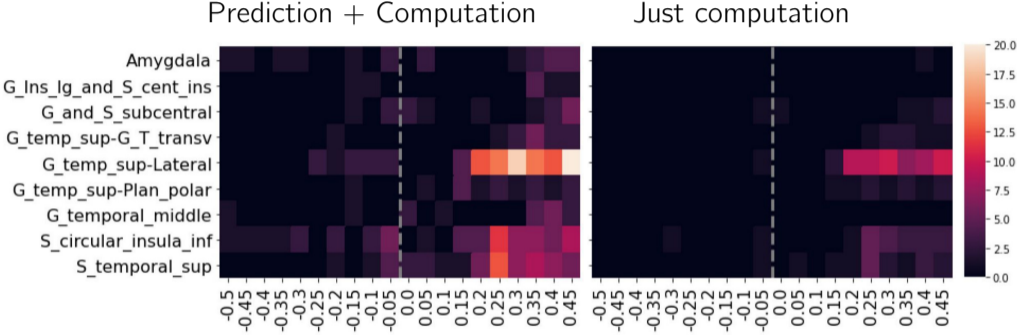


# Prediction vs computation



Balance data to make previous POS not predictive about next POS

# Prediction vs computation



Balance data to make previous POS not predictive about next POS  
Areas that predict POS are now gone, isolating POS computation

# The long road ahead . . .



# The long road ahead . . .

Coherent stories

# The long road ahead . . .

Coherent stories

3D

# The long road ahead . . .

Coherent stories

3D

Physics: Forces & contact relations



# Physics



# Physics



# The long road ahead . . .

Coherent stories

3D

Physics: Forces & contact relations

# The long road ahead . . .

Coherent stories

3D

Physics: Forces & contact relations

Segmentation



# The long road ahead . . .

Coherent stories

3D

Physics: Forces & contact relations

Segmentation

Parts and low-level features



# The long road ahead . . .

Coherent stories

3D

Physics: Forces & contact relations

Segmentation

Parts and low-level features

# The long road ahead . . .

Coherent stories

3D

Physics: Forces & contact relations

Segmentation

Parts and low-level features

Theory of mind



# The long road ahead . . .

Coherent stories

3D

Physics: Forces & contact relations

Segmentation

Parts and low-level features

Theory of mind

# The long road ahead . . .

Coherent stories

3D

Physics: Forces & contact relations

Segmentation

Parts and low-level features

Theory of mind

Social understanding

# Theory of mind





# Theory of mind



# The long road ahead . . .

Coherent stories

3D

Physics: Forces & contact relations

Segmentation

Parts and low-level features

Theory of mind

Social understanding

# The long road ahead . . .

Coherent stories

3D

Physics: Forces & contact relations

Segmentation

Parts and low-level features

Theory of mind

Social understanding

Modification







Boris Katz



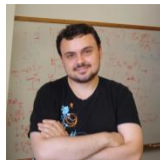
Shimon Ullman



Josh Tenenbaum



Gabriel Kreiman



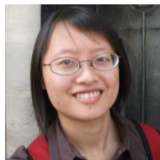
Andrei Barbu



Ignacio Cases



Candace Ross



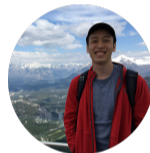
Yen-Ling Kuo



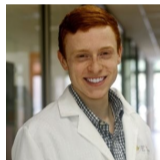
Adam Yaari



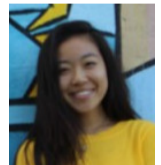
David Mayo



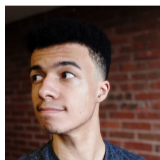
Christopher Wang



Julian Alverio



Emily Cheng



Dylan Sleeper



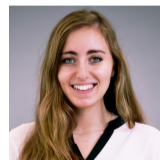
Vighnesh Subramaniam



Aaditya Singh



Ravi Tejwani



Dana Rosenfarb

Recognition	$P(\text{sentence}, \text{video})$	Narayanaswamy et al. 2014
Retrieval	$\operatorname{argmax}_{v \in V} P(s, v)$	Barret et al. 2016
Generation	$\operatorname{argmax}_{s \in L} P(s, v)$	Yu et al. 2015, N. et al. 2014
Question answering	$\operatorname{argmax}_{s \in L} P(Q(s, s_q), v)$	Barbu et al. in prep.
Disambiguation	$\operatorname{argmax}_{i \in \text{parser}(s)} P(i, v)$	Berzak et al. 2015
Language acquisition	$\operatorname{argmax}_{\theta} \prod_{s, v} P(s(\theta), v)$	Yu et al. 2015, Ross et al. 2018
Paraphrasing	$\int_v  P(s, v) - P(s', v) $	Mao et al. in review
Translation	$\operatorname{argmin}_{s' \in L'} \int_v  P(s, v) - P(s', v) $	Fu et al. in prep.
Common sense reasoning	$\operatorname{argmax}_{s \in L} \int_v P(s_q, v) P(Q(s, s_q), v)$	
Planning	$\operatorname{argmax}_{s \in L} \int_v P(s, v_0 : v : v_n)$	Kuo et al. 2018, Kuo et al. 2020
Command following	$\operatorname{argmax}_p \int_{v^+} P(C(s), v^+ v) E(v^+, p, v)$	Paul et al. 2017, Kuo et al. in prep.

There are deep connections between language, vision, perception, and embodiment



There are deep connections between language, vision, perception, and embodiment  
Most research is on supervised uni-modal uni-task in-domain models.

There are deep connections between language, vision, perception, and embodiment  
Most research is on supervised uni-modal uni-task in-domain models.  
But our brains are unlikely to be anything like those models!

There are deep connections between language, vision, perception, and embodiment

Most research is on supervised uni-modal uni-task in-domain models.

But our brains are unlikely to be anything like those models!

We need to take seriously zero-shot multi-modal multi-task out-of-domain learning.

There are deep connections between language, vision, perception, and embodiment

Most research is on supervised uni-modal uni-task in-domain models.

But our brains are unlikely to be anything like those models!

We need to take seriously zero-shot multi-modal multi-task out-of-domain learning.

The dataset crisis in vision & NLP is just a manifestation of this.

There are deep connections between language, vision, perception, and embodiment

Most research is on supervised uni-modal uni-task in-domain models.

But our brains are unlikely to be anything like those models!

We need to take seriously zero-shot multi-modal multi-task out-of-domain learning.

The dataset crisis in vision & NLP is just a manifestation of this.

These domains are all compositional and that looks to be our biggest hammer.

There are deep connections between language, vision, perception, and embodiment

Most research is on supervised uni-modal uni-task in-domain models.

But our brains are unlikely to be anything like those models!

We need to take seriously zero-shot multi-modal multi-task out-of-domain learning.

The dataset crisis in vision & NLP is just a manifestation of this.

These domains are all compositional and that looks to be our biggest hammer.

But we have no idea about why and how compositionality works.

There are deep connections between language, vision, perception, and embodiment

Most research is on supervised uni-modal uni-task in-domain models.

But our brains are unlikely to be anything like those models!

We need to take seriously zero-shot multi-modal multi-task out-of-domain learning.

The dataset crisis in vision & NLP is just a manifestation of this.

These domains are all compositional and that looks to be our biggest hammer.

But we have no idea about why and how compositionality works.

We use many crutches (like data augmentation) to avoid dealing with the theory of compositionality.

There are deep connections between language, vision, perception, and embodiment

Most research is on supervised uni-modal uni-task in-domain models.

But our brains are unlikely to be anything like those models!

We need to take seriously zero-shot multi-modal multi-task out-of-domain learning.

The dataset crisis in vision & NLP is just a manifestation of this.

These domains are all compositional and that looks to be our biggest hammer.

But we have no idea about why and how compositionality works.

We use many crutches (like data augmentation) to avoid dealing with the theory of compositionality.

Maybe flexible intelligence is possible because it's multi-modal and multi-task?



There are deep connections between language, vision, perception, and embodiment

Most research is on supervised uni-modal uni-task in-domain models.

But our brains are unlikely to be anything like those models!

We need to take seriously zero-shot multi-modal multi-task out-of-domain learning.

The dataset crisis in vision & NLP is just a manifestation of this.

These domains are all compositional and that looks to be our biggest hammer.

But we have no idea about why and how compositionality works.

We use many crutches (like data augmentation) to avoid dealing with the theory of compositionality.

Maybe flexible intelligence is possible because it's multi-modal and multi-task? Maybe our poor understanding of the brain is a consequence of building narrow models?