# Understanding vision through language and language through vision



Andrei Barbu

# Language and perception

# Language and perception



Caption, answer questions, understand a description, explain it to someone, engage in a conversation, give agents commands, imagine something different, recognize its description in a story, rewrite that description in another language, understand if someone is missing the point, reproduce it, intervene, etc.

Recognition

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

Follow commands

Recognition

Retrieval

Generation

Question answering

Disambiguation

Language acquisition

Follow commands

Paraphrasing

Recognition

Retrieval

Generation

Question answering

Disambiguation

Language acquisition

Follow commands

Paraphrasing

Translation

Recognition

Retrieval

Generation

Question answering

Disambiguation

Language acquisition

Follow commands

Paraphrasing

Translation

Common sense reasoning

Recognition

Retrieval

Generation

Question answering

Disambiguation

Language acquisition

Follow commands

Paraphrasing

Translation

Common sense reasoning

Planning

Recognition

Retrieval

Generation

Question answering

Disambiguation

Language acquisition

Follow commands

Paraphrasing

Translation

Common sense reasoning

Planning

  …

Recognition

Retrieval

Generation

Question answering                Computer vision

Disambiguation                         NLP

                                              Robotics
Language acquisition                    AI

Follow commands

Paraphrasing

Translation

Common sense reasoning

Planning

  ...

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

Narayanaswamy *et al.* 2014

Retrieval

Generation

Question answering

Disambiguation

Language acquisition

Follow commands

Paraphrasing

Translation

Common sense reasoning

Planning

  ...

The person rode the skateboard leftward.

The person rode the skateboard leftward.

object detector, tracker, event recognizer

The person rode the skateboard leftward.

object detector, tracker, event recognizer

The person rode the skateboard leftward.
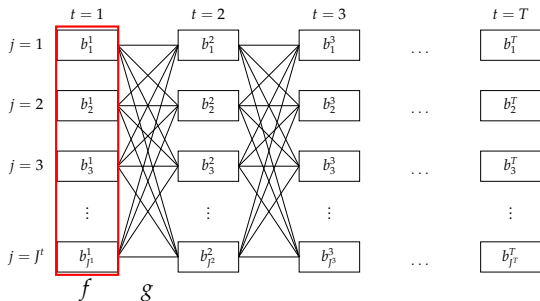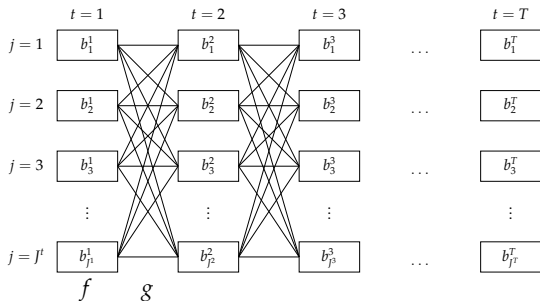
object detector, tracker, event recognizer

The person rode the skateboard leftward.

object detector, tracker, event recognizer

# Tracking with higher-level knowledge
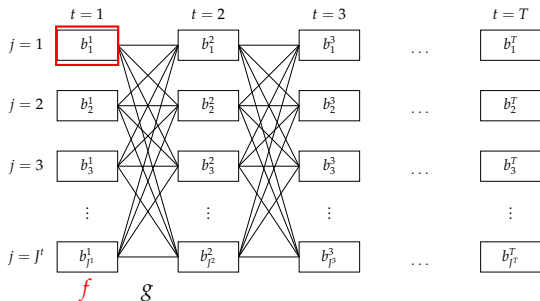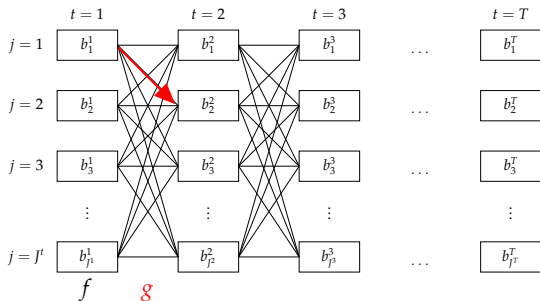
# Tracking with higher-level knowledge



detection / object / frame

Barbu *et al.* 2012

# Tracking with higher-level knowledge



detection / object / frame

temporally coherent track

Barbu *et al.* 2012
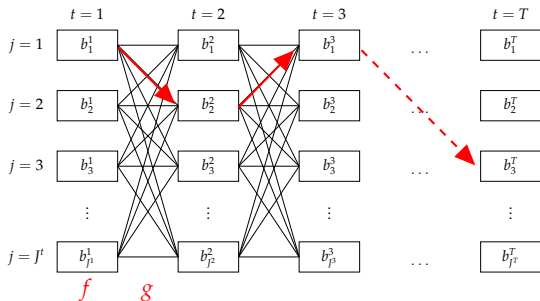
# Tracking with higher-level knowledge



detection / object / frame

temporally coherent track

object detector confidence (f)

Barbu *et al.* 2012

# Tracking with higher-level knowledge



detection / object / frame

temporally coherent track

    object detector confidence (f)

    motion coherence (g)

Barbu *et al.* 2012

# Tracking with higher-level knowledge



detection / object / frame

temporally coherent track

    object detector confidence (f)

    motion coherence (g)

optimal path through the lattice of detections

Barbu *et al.* 2012

# Tracking with higher-level knowledge



detection / object / frame

temporally coherent track
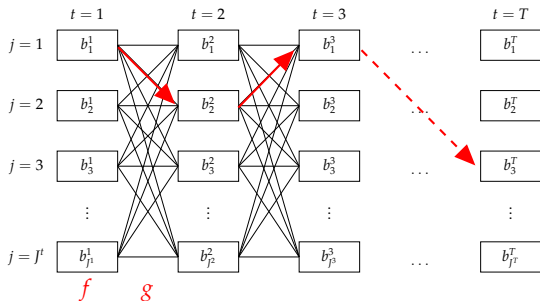
> object detector confidence (f)
>
> motion coherence (g)

optimal path through the lattice of detections

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

Barbu *et al.* 2012

# Tracking with higher-level knowledge



detection / object / frame

temporally coherent track

- object detector confidence (f)
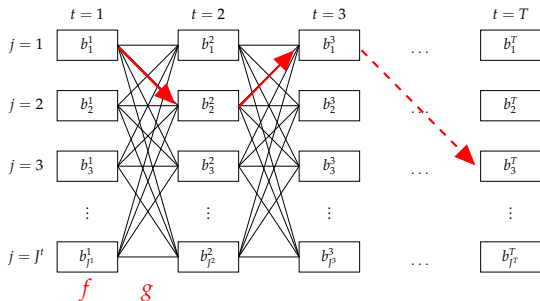- motion coherence (g)

optimal path through the lattice of detections

dynamic programming

Bellman (1957), Viterbi (1967)

$$\max_{j^1,\ldots,j^T} \sum_{t=1}^{T} f(b_{j^t}^t) + \sum_{t=2}^{T} g(b_{j^{t-1}}^{t-1}, b_{j^t}^t)$$
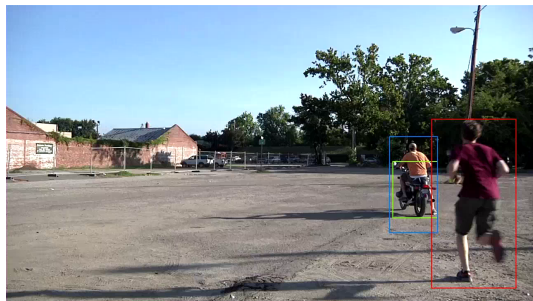
Barbu *et al.* 2012

# Tracking 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)

$$\max_{j^1,\ldots,j^T} \sum_{t=1}^{T} f(b_{j^t}^t) + \sum_{t=2}^{T} g(b_{j^{t-1}}^{t-1}, b_{j^t}^t)$$
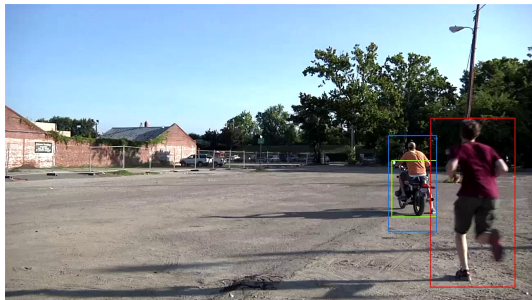
Barbu *et al.* 2012

# Tracking 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)

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

Barbu *et al.* 2012

# Event recognition

 tracks

# Event recognition



feature vectors

tracks

# Event recognition



$a$

HMMs

Baum and Petrie (1966)

feature vectors

tracks

# Event recognition



*a*  HMMs

Baum and Petrie (1966)

*h*

feature vectors

tracks

# Event recognition



*a*    HMMs

Baum and Petrie (1966)

*h*

feature vectors

tracks

# Event recognition



HMMs
Baum and Petrie (1966)
One HMM per event class

feature vectors

tracks

# Event recognition



*a* HMMs
Baum and Petrie (1966)
One HMM per event class
Try each HMM

*h*

feature vectors

tracks

# Event recognition

# Event recognition

# Event recognition



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

# Building sentences out of trackers and words

Siddharth *et al.* 2014

# Building sentences out of trackers and words

Viterbi tracker

track 1



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

Siddharth *et al.* 2014

# Building sentences out of trackers and words

Event tracker for intransitive verbs

track 1



$\times$



word 1

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

Siddharth *et al.* 2014

# Building sentences out of trackers and words



Event tracker

track 1         track $L$

$\times \cdots \times$

$\times$

word 1

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

Siddharth *et al.* 2014

# Building sentences out of trackers and words

Sentence tracker



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

Siddharth *et al.* 2014

# Sentence recognizers

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

# Sentence recognizers

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



. . .

# Sentence recognizers

The tall <span style="color:red">person</span> quickly rode the <span style="color:red">horse</span> leftward away from the other <span style="color:red">horse.</span>



. . .

# Sentence recognizers

The tall <span style="color:red">person</span> quickly rode the <span style="color:red">horse</span> leftward away from the other <span style="color:red">horse.</span>



agent-location    patient-location

...

# Sentence recognizers

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

agent-location    patient-location



. . .

# Sentence recognizers

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

agent-location        patient-location



...

# Sentence recognizers

The tall person quickly rode the <span style="color:red">horse</span> leftward away from the <span style="color:blue">other</span> <span style="color:red">horse</span>.

agent-location    patient-location



. . .

# Sentence recognizers

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

agent-location    patient-location    source-location



...

# Sentence recognizers

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

agent-location    patient-location    source-location



...

# Sentence recognizers



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

agent-location    patient-location    source-location

...

# Sentence recognizers



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

agent-location    patient-location    source-location

# Sentence recognizers



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

agent-location    patient-location    source-location

. . .

# Sentence recognizers

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

agent-location        patient-location        source-location



...

# Sentence recognizers



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

agent-location    patient-location    source-location

...

# Sentence recognizers



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

agent-location    patient-location    source-location

...

The person picked up the animal from the table.

The person picked up the animal from the table.

$\exists xyz\ \textbf{person}(x), \textbf{animal}(y), \textbf{table}(z), \textbf{pickup}(x, y), \textbf{on}(y, z)$

The person picked up the animal from the table.

$\exists xyz\ \textbf{person}(x), \textbf{animal}(y), \textbf{table}(z), \textbf{pickup}(x, y), \textbf{on}(y, z)$

**chair**    **animal**    **table**    **pickup**    **on**

$v_{xy}, a_{xy}, \ldots$        $v_{xy}, a_{xy}, \ldots$        $v_{xy}, a_{xy}, \ldots$

$x$-tracker            $y$-tracker            $z$-tracker

The person carried something.
The person went away.
The person walked.
The person had the bag.
The person left leftward and upward.

Barbu *et al.* 2012

The person carried something.
The person went away.
The person walked.
The person had the bag.
The person left leftward and upward.

Barbu *et al.* 2012

Retrieval

Generation

Question answering

Disambiguation

Language acquisition

Follow commands

Paraphrasing

Translation

Common sense reasoning

Planning

  …

| | | |
|---|---|---|
| Recognition | $P(\text{sentence}, \text{video})$ | Narayanaswamy *et al.* 2014 |
| Retrieval | $\underset{v \in V}{\arg\max}\; P(s, v)$ | Barret *et al.* 2016 |
| Generation | | |
| Question answering | | |
| Disambiguation | | |
| Language acquisition | | |
| Follow commands | | |
| Paraphrasing | | |
| Translation | | |
| Common sense reasoning | | |
| Planning | | |
| … | | |

# Sentential retrieval

# Sentential retrieval

Recognition          $P(\text{sentence}, \text{video})$                    Narayanaswamy *et al.* 2014

Retrieval            $\underset{v \in V}{\text{argmax}} \; P(s, v)$          Barret *et al.* 2016

Generation

Question answering

Disambiguation

Language acquisition

Follow commands

Paraphrasing

Translation

Common sense reasoning

Planning

    …

| Recognition | $P(\text{sentence}, \text{video})$ | Narayanaswamy *et al.* 2014 |
| Retrieval | $\underset{v \in V}{\arg\max}\ P(s, v)$ | Barret *et al.* 2016 |
| Generation | $\underset{s \in L}{\arg\max}\ P(s, v)$ | Yu *et al.* 2015, Narayanaswamy *et al.* 2 |
| Question answering | | |
| Disambiguation | | |
| Language acquisition | | |
| Follow commands | | |
| Paraphrasing | | |
| Translation | | |
| Common sense reasoning | | |
| Planning | | |
| … | | |

# Generating sentences

# Generating sentences

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

# Generating sentences

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

147,123,874,800 sentences without recursion

# Generating sentences

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

147,123,874,800 sentences without recursion

$\emptyset$

# Generating sentences

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

147,123,874,800 sentences without recursion

"carried"

# Generating sentences

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

147,123,874,800 sentences without recursion

"the person carried"

# Generating sentences

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

147,123,874,800 sentences without recursion

"the person carried the backpack"

# Generated sentences

# Generated sentences

# Generated sentences



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

| | | |
|---|---|---|
| Recognition | $P(\text{sentence}, \text{video})$ | Narayanaswamy *et al.* 2014 |
| Retrieval | $\underset{v \in V}{\text{argmax}} \; P(s, v)$ | Barret *et al.* 2016 |
| Generation | $\underset{s \in L}{\text{argmax}} \; P(s, v)$ | Yu *et al.* 2015, Narayanaswamy *et al.* 2 |
| Question answering | | |
| Disambiguation | | |
| Language acquisition | | |
| Follow commands | | |
| Paraphrasing | | |
| Translation | | |
| Common sense reasoning | | |
| Planning | | |
| … | | |

| Recognition | $P(\text{sentence}, \text{video})$ | Narayanaswamy *et al.* 2014 |
| Retrieval | $\underset{v \in V}{\text{argmax}}\ P(s, v)$ | Barret *et al.* 2016 |
| Generation | $\underset{s \in L}{\text{argmax}}\ P(s, v)$ | Yu *et al.* 2015, Narayanaswamy *et al.* 2 |
| Question answering | $\underset{s \in L}{\text{argmax}}\ P(\text{Q}(s, s_q), v)$ | Barbu *et al.* in prep. |
| Disambiguation | | |
| Language acquisition | | |
| Follow commands | | |
| Paraphrasing | | |
| Translation | | |
| Common sense reasoning | | |
| Planning | | |
| ... | | |

# 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.

# Question answering

# Question answering

# Question answering

# Question answering



Who put an object on top of the red car?

# Question answering with specificity



Who put an object on top of the red car?
NP put an object on top of the red car.

# Question answering with specificity



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 | $P(\text{sentence}, \text{video})$ | Narayanaswamy *et al.* 2014 |
| Retrieval | $\underset{v \in V}{\text{argmax}}\ P(s, v)$ | Barret *et al.* 2016 |
| Generation | $\underset{s \in L}{\text{argmax}}\ P(s, v)$ | Yu *et al.* 2015, Narayanaswamy *et al.* 2 |
| Question answering | $\underset{s \in L}{\text{argmax}}\ P(\text{Q}(s, s_q), v)$ | Barbu *et al.* in prep. |
| Disambiguation | | |
| Language acquisition | | |
| Follow commands | | |
| Paraphrasing | | |
| Translation | | |
| Common sense reasoning | | |
| Planning | | |
| ... | | |

| | | |
|---|---|---|
| Recognition | $P(\text{sentence}, \text{video})$ | Narayanaswamy *et al.* 2014 |
| Retrieval | $\underset{v \in V}{\text{argmax}}\ P(s, v)$ | Barret *et al.* 2016 |
| Generation | $\underset{s \in L}{\text{argmax}}\ P(s, v)$ | Yu *et al.* 2015, Narayanaswamy *et al.* 2 |
| Question answering | $\underset{s \in L}{\text{argmax}}\ P(\text{Q}(s, s_q), v)$ | Barbu *et al.* in prep. |
| Disambiguation | $\underset{i \in \text{parser}(s)}{\text{argmax}}\ P(i, v)$ | Berzak *et al.* 2015 |
| Language acquisition | | |
| Follow commands | | |
| Paraphrasing | | |
| Translation | | |
| Common sense reasoning | | |
| Planning | | |
| … | | |

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
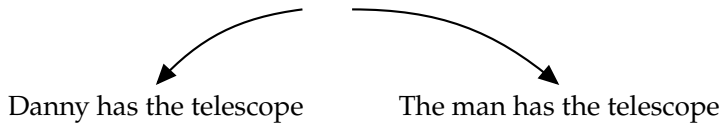
Danny looked at the man with a telescope.

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

VP Attachment

Andrei approached the person holding a green chair.

# 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.

Danny and Andrei moved a chair.

Danny and Andrei moved a chair.

Danny and Andrei move the same chair.

Danny and Andrei moved a chair.

Danny and Andrei move the same chair.
$$\exists x \; \mathbf{chair}(x)$$
$$\mathbf{move}(\text{Danny}, x), \mathbf{move}(\text{Andrei}, x)$$

# Danny and Andrei moved a chair.

Danny and Andrei move the same chair.
$$\exists x \; \textbf{chair}(x)$$
$$\textbf{move}(\text{Danny}, x), \textbf{move}(\text{Andrei}, x)$$

Danny and Andrei move different chairs.

# Danny and Andrei moved a chair.

Danny and Andrei move the same chair.
$$\exists x \ \mathbf{chair}(x)$$
$$\mathbf{move}(\text{Danny}, x), \mathbf{move}(\text{Andrei}, x)$$

Danny and Andrei move different chairs.
$$\exists xy \ \mathbf{chair}(x), \mathbf{chair}(y)$$
$$\mathbf{move}(\text{Danny}, x), \mathbf{move}(\text{Andrei}, y), x \neq y$$

| | | |
|---|---|---|
| Recognition | $P(\text{sentence}, \text{video})$ | Narayanaswamy *et al.* 2014 |
| Retrieval | $\underset{v \in V}{\text{argmax}} \ P(s, v)$ | Barret *et al.* 2016 |
| Generation | $\underset{s \in L}{\text{argmax}} \ P(s, v)$ | Yu *et al.* 2015, Narayanaswamy *et al.* 2 |
| Question answering | $\underset{s \in L}{\text{argmax}} \ P(\text{Q}(s, s_q), v)$ | Barbu *et al.* in prep. |
| Disambiguation | $\underset{i \in \text{parser}(s)}{\text{argmax}} \ P(i, v)$ | Berzak *et al.* 2015 |
| Language acquisition | | |
| Follow commands | | |
| Paraphrasing | | |
| Translation | | |
| Common sense reasoning | | |
| Planning | | |
| … | | |

| Recognition | $P(\text{sentence}, \text{video})$ | Narayanaswamy *et al.* 2014 |
| Retrieval | $\underset{v \in V}{\text{argmax}}\ P(s, v)$ | Barret *et al.* 2016 |
| Generation | $\underset{s \in L}{\text{argmax}}\ P(s, v)$ | Yu *et al.* 2015, Narayanaswamy *et al.* 2 |
| Question answering | $\underset{s \in L}{\text{argmax}}\ P(Q(s, s_q), v)$ | Barbu *et al.* in prep. |
| Disambiguation | $\underset{i \in \text{parser}(s)}{\text{argmax}}\ P(i, v)$ | Berzak *et al.* 2015 |
| Language acquisition | $\underset{\theta}{\text{argmax}} \prod_{s,v} P(s(\theta), v)$ | Yu *et al.* 2015 |

Follow commands

Paraphrasing

Translation

Common sense reasoning

Planning

…

# Language learning

# Language learning

Split into two variants:

# Language learning

Split into two variants:

    Lexicon

# Language learning

Split into two variants:
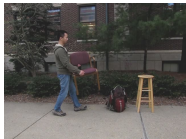Lexicon
Syntax

# Language learning

Split into two variants:

Lexicon
Syntax

# Language learning: Lexicon

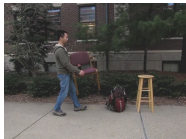Yu *et al.* 2014

# Language learning: Lexicon



Yu *et al.* 2014

# Language learning: Lexicon



The person picked up the chair.



The chair approached the backpack.



The person picked up the backpack.

Yu *et al.* 2014

# Language learning: Lexicon



The person picked up the chair.



The chair approached the backpack.



The person picked up the backpack.

Yu *et al.* 2014

# Language learning: Lexicon



The person picked up the chair.



The chair approached the backpack.



The person picked up the backpack.

Yu *et al.* 2014

# Language learning: Lexicon



chair

The person picked up the chair.

picked up

The chair approached the backpack.

person

approached

The person picked up the backpack.

backpack

Yu *et al.* 2014

# Language learning

Split into two variants:

Lexicon
Syntax

# Language learning

Split into two variants:

Lexicon

Syntax

# Language learning: Syntax

# Language learning: Syntax

Danny approached the chair with a bag.

# Language learning: Syntax

Danny approached the chair with a bag.

parser

# Language learning: Syntax

Danny approached the chair with a bag.



parser

# Language learning: Syntax

Danny approached the chair with a bag.



parser

Danny approached the chair with a bag.

parser

# Language learning: Syntax



Danny approached the chair with a bag.

parser

# Language learning: Syntax



Danny approached the chair with a bag.

$\approx$parser

$\bullet \bullet \bullet$

# Language learning: Syntax



Danny approached the chair with a bag.
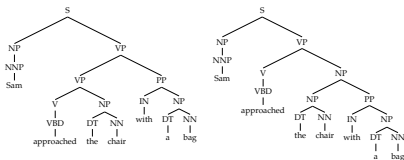
$\approx$parser

# Language learning: Syntax

Danny approached the chair with a bag.



≈parser



• • •

# Language learning: Syntax



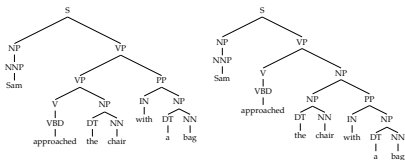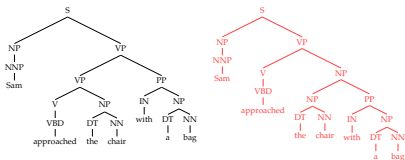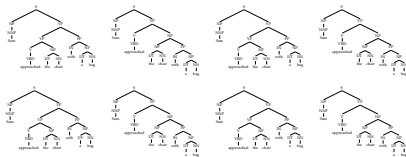Danny approached the chair with a bag.

$\downarrow$

$\approx$parser

$\downarrow$

Pilley and Reid 2011

Pilley and Reid 2011

| | | |
|---|---|---|
| Recognition | $P(\text{sentence}, \text{video})$ | Narayanaswamy *et al.* 2014 |
| Retrieval | $\underset{v \in V}{\text{argmax}}\ P(s, v)$ | Barret *et al.* 2016 |
| Generation | $\underset{s \in L}{\text{argmax}}\ P(s, v)$ | Yu *et al.* 2015, Narayanaswamy *et al.* 2 |
| Question answering | $\underset{s \in L}{\text{argmax}}\ P(\text{Q}(s, s_q), v)$ | Barbu *et al.* in prep. |
| Disambiguation | $\underset{i \in \text{parser}(s)}{\text{argmax}}\ P(i, v)$ | Berzak *et al.* 2015 |
| Language acquisition | $\underset{\theta}{\text{argmax}} \prod_{s,v} P(s(\theta), v)$ | Yu *et al.* 2015 |
| Follow commands | | |
| Paraphrasing | | |
| Translation | | |
| Common sense reasoning | | |
| Planning | | |
| … | | |

| | | |
|---|---|---|
| Recognition | $P(\text{sentence}, \text{video})$ | Narayanaswamy *et al.* 2014 |
| Retrieval | $\underset{v \in V}{\text{argmax}} \ P(s, v)$ | Barret *et al.* 2016 |
| Generation | $\underset{s \in L}{\text{argmax}} \ P(s, v)$ | Yu *et al.* 2015, Narayanaswamy *et al.* 2 |
| Question answering | $\underset{s \in L}{\text{argmax}} \ P(\text{Q}(s, s_q), v)$ | Barbu *et al.* in prep. |
| Disambiguation | $\underset{i \in \text{parser}(s)}{\text{argmax}} \ P(i, v)$ | Berzak *et al.* 2015 |
| Language acquisition | $\underset{\theta}{\text{argmax}} \prod_{s, v} P(s(\theta), v)$ | Yu *et al.* 2015 |
| Follow commands | $\underset{p}{\text{argmax}} \int_{v^+} P(C(s), v^+ v) \ E(v^+, p, v)$ | Paul *et al.* 2017 |
| Paraphrasing | | |
| Translation | | |
| Common sense reasoning | | |
| Planning | | |
| … | | |

... ...

The box I put down on the table is expensive.

 ...  ...

The box I put down on the table is expensive.

$\text{Human}(x) \wedge \text{PutDown}(x, y) \wedge \text{Box}(y) \wedge \text{On}(y, z) \wedge \text{Table}(z)$      $\text{Assert}(\text{Expensive}(y))$

The box I put down on the table is expensive.

$Human(x) \wedge PutDown(x, y) \wedge Box(y) \wedge On(y, z) \wedge Table(z)$    $Assert(Expensive(y))$

The box I put down on the table is expensive.

Human($x$) $\wedge$ PutDown($x, y$) $\wedge$ Box($y$) $\wedge$ On($y, z$) $\wedge$ Table($z$)          Assert(Expensive($y$))



⊥
State

The box I put down on the table is expensive.

$$\text{Human}(x) \wedge \text{PutDown}(x, y) \wedge \text{Box}(y) \wedge \text{On}(y, z) \wedge \text{Table}(z)$$

$$\text{Assert}(\text{Expensive}(y))$$



⊥

State

Pick up the expensive box someone put down next to the can.

The box I put down on the table is expensive.

Human($x$) ∧ PutDown($x, y$) ∧ Box($y$) ∧ On($y, z$) ∧ Table($z$)     Assert(Expensive($y$))



State

Pick up the expensive box someone put down next to the can.



...

The box I put down on the table is expensive.

$\text{Human}(x) \wedge \text{PutDown}(x, y) \wedge \text{Box}(y) \wedge \text{On}(y, z) \wedge \text{Table}(z)$

$\text{Assert}(\text{Expensive}(y))$



⊥

State

Pick up the expensive box someone put down next to the can.

$\text{PutDown}(x, y) \wedge \text{Expensive}(y) \wedge \text{Can}(z) \wedge \text{Box}(y) \wedge \text{NextTo}(y, z)$

The box I put down on the table is expensive.

$Human(x) \land PutDown(x, y) \land Box(y) \land On(y, z) \land Table(z)$    $Assert(Expensive(y))$



$\bot$
State

Pick up the expensive box someone put down next to the can.

$PutDown(x, y) \land Expensive(y) \land Can(z) \land Box(y) \land NextTo(y, z)$
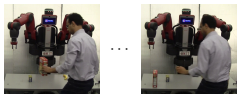
The box I put down on the table is expensive.

$\text{Human}(x) \wedge \text{PutDown}(x, y) \wedge \text{Box}(y) \wedge \text{On}(y, z) \wedge \text{Table}(z)$

$\text{Assert}(\text{Expensive}(y))$



⊥

State

Pick up the expensive box someone put down next to the can.

$\text{PutDown}(x, y) \wedge \text{Expensive}(y) \wedge \text{Can}(z) \wedge \text{Box}(y) \wedge \text{NextTo}(y, z)$

The box I put down on the table is expensive.

Human($x$) ∧ PutDown($x, y$) ∧ Box($y$) ∧ On($y, z$) ∧ Table($z$)     Assert(Expensive($y$))



State

Pick up the expensive box someone put down next to the can.

PutDown($x, y$) ∧ Expensive($y$) ∧ Can($z$) ∧ Box($y$) ∧ NextTo($y, z$)



$\lambda_1$          $\lambda_2$
Pick up          box

The box I put down on the table is expensive.

$Human(x) \wedge PutDown(x, y) \wedge Box(y) \wedge On(y, z) \wedge Table(z)$          $Assert(Expensive(y))$



State

Pick up the expensive box someone put down next to the can.

$PutDown(x, y) \wedge Expensive(y) \wedge Can(z) \wedge Box(y) \wedge NextTo(y, z)$



$\gamma_1$   $\gamma_2$

$\phi_1$   $\phi_2$

$\lambda_1$   $\lambda_2$

Pick up   box

The box I put down on the table is expensive.

$\text{Human}(x) \wedge \text{PutDown}(x, y) \wedge \text{Box}(y) \wedge \text{On}(y, z) \wedge \text{Table}(z)$          $\text{Assert}(\text{Expensive}(y))$
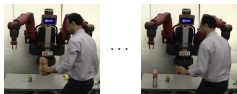
State

Pick up the expensive box someone put down next to the can.

$\text{PutDown}(x, y) \wedge \text{Expensive}(y) \wedge \text{Can}(z) \wedge \text{Box}(y) \wedge \text{NextTo}(y, z)$

$\gamma_1$          $\gamma_2$

$\phi_1$          $\phi_2$

$\lambda_1$          $\lambda_2$

Pick up          box

The box I put down on the table is expensive.

$\text{Human}(x) \wedge \text{PutDown}(x, y) \wedge \text{Box}(y) \wedge \text{On}(y, z) \wedge \text{Table}(z)$      $\text{Assert}(\text{Expensive}(y))$
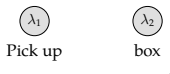
State

Pick up the expensive box someone put down next to the can.

$\text{PutDown}(x, y) \wedge \text{Expensive}(y) \wedge \text{Can}(z) \wedge \text{Box}(y) \wedge \text{NextTo}(y, z)$

$\gamma_1$     $\gamma_2$

$\phi_1$     $\phi_2$

$\lambda_1$     $\lambda_2$

Pick up     box

# Understanding and following commands

# Understanding and following commands



Semantic parser

# Understanding and following commands

# Understanding and following commands



Semantic parser

Perception

# Understanding and following commands

# Understanding and following commands

The speaker informs the robot,
"the box I will put down is my snack".

The speaker then places a box on the table.

The speaker informs the robot,
"the box I will put down is my snack".

The speaker then places a box on the table.

The speaker informs the robot,
"the box I will put down is my snack".

The speaker then places a box on the table.

| | | |
|---|---|---|
| Recognition | $P(\text{sentence}, \text{video})$ | Narayanaswamy *et al.* 2014 |
| Retrieval | $\underset{v \in V}{\text{argmax}} \; P(s, v)$ | Barret *et al.* 2016 |
| Generation | $\underset{s \in L}{\text{argmax}} \; P(s, v)$ | Yu *et al.* 2015, Narayanaswamy *et al.* 2 |
| Question answering | $\underset{s \in L}{\text{argmax}} \; P(\text{Q}(s, s_q), v)$ | Barbu *et al.* in prep. |
| Disambiguation | $\underset{i \in \text{parser}(s)}{\text{argmax}} \; P(i, v)$ | Berzak *et al.* 2015 |
| Language acquisition | $\underset{\theta}{\text{argmax}} \prod_{s, v} P(s(\theta), v)$ | Yu *et al.* 2015 |
| Follow commands | $\underset{p}{\text{argmax}} \int_{v^+} P(C(s), v^+ v) \; E(v^+, p, v)$ | Paul *et al.* 2017 |
| Paraphrasing | | |
| Translation | | |
| Common sense reasoning | | |
| Planning | | |
| … | | |

| | | |
|---|---|---|
| Recognition | $P(\text{sentence}, \text{video})$ | Narayanaswamy *et al.* 2014 |
| Retrieval | $\underset{v \in V}{\arg\max}\ P(s, v)$ | Barret *et al.* 2016 |
| Generation | $\underset{s \in L}{\arg\max}\ P(s, v)$ | Yu *et al.* 2015, Narayanaswamy *et al.* 2 |
| Question answering | $\underset{s \in L}{\arg\max}\ P(\text{Q}(s, s_q), v)$ | Barbu *et al.* in prep. |
| Disambiguation | $\underset{i \in \text{parser}(s)}{\arg\max}\ P(i, v)$ | Berzak *et al.* 2015 |
| Language acquisition | $\underset{\theta}{\arg\max} \prod_{s,v} P(s(\theta), v)$ | Yu *et al.* 2015 |
| Follow commands | $\underset{p}{\arg\max} \int_{v^+} P(C(s), v^+ v)\, E(v^+, p, v)$ | Paul *et al.* 2017 |
| Paraphrasing | $\int_v \lvert\, P(s, v)\, - \, P(s', v)\, \rvert$ | |
| Translation | | |
| Common sense reasoning | | |
| Planning | | |
| ... | | |

# Paraphrasing

# Paraphrasing

The dark haired man is picking an object up from the floor.
The guy in the plaid shirt is lifting the yellow chair.

# Paraphrasing

The dark haired man is picking an object up from the floor.
The guy in the plaid shirt is lifting the yellow chair.


The man with the chair walks away from someone.
The man walks away from someone with the chair.

# Paraphrasing

# Paraphrasing

The tall man gave the woman the red box.

# Paraphrasing

The tall man gave the woman the red box.

$\overset{?}{\Rightarrow}$ The woman received the crimson box from the man.

# Paraphrasing

The tall man gave the woman the red box.

$\overset{?}{\Rightarrow}$ The woman received the crimson box from the man.

$\overset{?}{\Rightarrow}$ The box was passed on by the body guard.

# Paraphrasing is hard

The tall man gave the woman the red box.

$\overset{?}{\Rightarrow}$ The woman received the crimson box from the man.

$\overset{?}{\Rightarrow}$ The box was passed on by the body guard.

$\overset{?}{\Rightarrow}$ The box changed hands from the man to the woman.

# Paraphrasing is hard

The tall man gave the woman the red box.

$\overset{?}{\Rightarrow}$ The woman received the crimson box from the man.

$\overset{?}{\Rightarrow}$ The box was passed on by the body guard.

$\overset{?}{\Rightarrow}$ The box changed hands from the man to the woman.

$\overset{?}{\Rightarrow}$ The man entered a room with a box. A woman left holding it.

# Paraphrasing is hard

The tall man gave the woman the red box.

$\overset{?}{\Rightarrow}$ The woman received the crimson box from the man.

$\overset{?}{\Rightarrow}$ The box was passed on by the body guard.

$\overset{?}{\Rightarrow}$ The box changed hands from the man to the woman.

$\overset{?}{\Rightarrow}$ The man entered a room with a box. A woman left holding it.

$\overset{?}{\Rightarrow}$ The cat yawned.

# Paraphrasing is hard

The tall man gave the woman the red box.

$\overset{?}{\Rightarrow}$ The woman received the crimson box from the man.

$\overset{?}{\Rightarrow}$ The box was passed on by the body guard.

$\overset{?}{\Rightarrow}$ The box changed hands from the man to the woman.

$\overset{?}{\Rightarrow}$ The man entered a room with a box. A woman left holding it.

$\overset{?}{\Rightarrow}$ The cat yawned.

$\overset{?}{\Rightarrow}$ The guard took the box from the woman.

# Paraphrasing is hard

The tall man gave the woman the red box.

$\overset{?}{\Rightarrow}$ The woman received the crimson box from the man.

$\overset{?}{\Rightarrow}$ The box was passed on by the body guard.

$\overset{?}{\Rightarrow}$ The box changed hands from the man to the woman.

$\overset{?}{\Rightarrow}$ The man entered a room with a box. A woman left holding it.

$\overset{?}{\Rightarrow}$ The cat yawned.

$\overset{?}{\Rightarrow}$ The guard took the box from the woman.

$\overset{?}{\Rightarrow}$ The woman left the box behind.

# Paraphrasing is hard

The tall man gave the woman the red box.

$\overset{?}{\Rightarrow}$ The woman received the crimson box from the man.

$\overset{?}{\Rightarrow}$ The box was passed on by the body guard.

$\overset{?}{\Rightarrow}$ The box changed hands from the man to the woman.

$\overset{?}{\Rightarrow}$ The man entered a room with a box. A woman left holding it.

$\overset{?}{\Rightarrow}$ The cat yawned.

$\overset{?}{\Rightarrow}$ The guard took the box from the woman.

$\overset{?}{\Rightarrow}$ The woman left the box behind.

$\overset{?}{\Rightarrow}$ The box was passed to the man.

# Paraphrasing is hard

The tall man gave the woman the red box.

$\overset{?}{\Rightarrow}$ The woman received the crimson box from the man.

$\overset{?}{\Rightarrow}$ The box was passed on by the body guard.

$\overset{?}{\Rightarrow}$ The box changed hands from the man to the woman.

$\overset{?}{\Rightarrow}$ The man entered a room with a box. A woman left holding it.

$\overset{?}{\Rightarrow}$ The cat yawned.

$\overset{?}{\Rightarrow}$ The guard took the box from the woman.

$\overset{?}{\Rightarrow}$ The woman left the box behind.

$\overset{?}{\Rightarrow}$ The box was passed to the man.

$\overset{?}{\Rightarrow}$ The crate was on the floor while the woman picked up the dog.

# Paraphrasing is hard

The tall man gave the woman the red box.

$\overset{?}{\Rightarrow}$ The woman received the crimson box from the man.

$\overset{?}{\Rightarrow}$ The box was passed on by the body guard.

$\overset{?}{\Rightarrow}$ The box changed hands from the man to the woman.

$\overset{?}{\Rightarrow}$ The man entered a room with a box. A woman left holding it.

$\overset{?}{\Rightarrow}$ The cat yawned.

$\overset{?}{\Rightarrow}$ The guard took the box from the woman.

$\overset{?}{\Rightarrow}$ The woman left the box behind.

$\overset{?}{\Rightarrow}$ The box was passed to the man.

$\overset{?}{\Rightarrow}$ The crate was on the floor while the woman picked up the dog.

$\overset{?}{\Rightarrow}$ The dog was on the floor while the woman picked up the crate.

# Paraphrasing is hard

The tall man gave the woman the red box.

$\overset{?}{\Rightarrow}$ The woman received the crimson box from the man.

$\overset{?}{\Rightarrow}$ The box was passed on by the body guard.

$\overset{?}{\Rightarrow}$ The box changed hands from the man to the woman.

$\overset{?}{\Rightarrow}$ The man entered a room with a box. A woman left holding it.

$\overset{?}{\Rightarrow}$ The cat yawned.

$\overset{?}{\Rightarrow}$ The guard took the box from the woman.

$\overset{?}{\Rightarrow}$ The woman left the box behind.

$\overset{?}{\Rightarrow}$ The box was passed to the man.

$\overset{?}{\Rightarrow}$ The crate was on the floor while the woman picked up the dog.

$\overset{?}{\Rightarrow}$ The dog was on the floor while the woman picked up the crate.

...

# Paraphrasing today

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

# Paraphrasing today



Paraphrase Softmax Classifier

Fixed-Sized Matrix

Dynamic Pooling Layer

Variable-Sized Similarity Matrix

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

# 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

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

YouTube

$s \longrightarrow$

# Paraphrasing with vision

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

YouTube



$s \longrightarrow$  $\longrightarrow$ Videos $v$

# Paraphrasing with vision

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

YouTube



$s \longrightarrow$     $\longrightarrow$ Videos $v$     $s'$

# Paraphrasing with vision

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



YouTube

$s \longrightarrow$ [videos] $\longrightarrow$ Videos $v$ $\qquad s'$

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

# Paraphrasing with vision

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



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

# Paraphrasing with vision

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



$$\sum_{v} |S(s,v) - S(s',v)|^2$$

# Paraphrasing with vision

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

$$s \xrightarrow{\text{sample}} \text{Videos } v \qquad s'$$

$$\underbrace{\qquad\qquad\qquad\qquad\qquad\qquad}_{\displaystyle \int_v |S(s,v) - S(s',v)|^2}$$

# Paraphrasing with imagination

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

$$s \xrightarrow{\text{sample}} \text{Videos } v \qquad s'$$

$$\underbrace{\qquad\qquad\qquad\qquad\qquad\qquad\qquad}_{\displaystyle \int_v |S(s,v) - S(s',v)|^2}$$

| | | |
|---|---|---|
| Recognition | $P(\text{sentence}, \text{video})$ | Narayanaswamy *et al.* 2014 |
| Retrieval | $\underset{v \in V}{\text{argmax}} \; P(s, v)$ | Barret *et al.* 2016 |
| Generation | $\underset{s \in L}{\text{argmax}} \; P(s, v)$ | Yu *et al.* 2015, Narayanaswamy *et al.* 2 |
| Question answering | $\underset{s \in L}{\text{argmax}} \; P(\text{Q}(s, s_q), v)$ | Barbu *et al.* in prep. |
| Disambiguation | $\underset{i \in \text{parser}(s)}{\text{argmax}} \; P(i, v)$ | Berzak *et al.* 2015 |
| Language acquisition | $\underset{\theta}{\text{argmax}} \prod_{s, v} P(s(\theta), v)$ | Yu *et al.* 2015 |
| Follow commands | $\underset{p}{\text{argmax}} \int_{v^+} P(C(s), v^+ v) \, E(v^+, p, v)$ | Paul *et al.* 2017 |
| Paraphrasing | $\int_v \mid P(s, v) - P(s', v) \mid$ | |
| Translation | | |
| Common sense reasoning | | |
| Planning | | |
| ... | | |

| | | |
|---|---|---|
| Recognition | $P(\text{sentence}, \text{video})$ | Narayanaswamy *et al.* 2014 |
| Retrieval | $\underset{v \in V}{\text{argmax}}\ P(s, v)$ | Barret *et al.* 2016 |
| Generation | $\underset{s \in L}{\text{argmax}}\ P(s, v)$ | Yu *et al.* 2015, Narayanaswamy *et al.* 2 |
| Question answering | $\underset{s \in L}{\text{argmax}}\ P(\text{Q}(s, s_q), v)$ | Barbu *et al.* in prep. |
| Disambiguation | $\underset{i \in \text{parser}(s)}{\text{argmax}}\ P(i, v)$ | Berzak *et al.* 2015 |
| Language acquisition | $\underset{\theta}{\text{argmax}} \prod_{s,v} P(s(\theta), v)$ | Yu *et al.* 2015 |
| Follow commands | $\underset{p}{\text{argmax}} \int_{v^+} P(C(s), v^+ v)\ E(v^+, p, v)$ | Paul *et al.* 2017 |
| Paraphrasing | $\int_v \mid P(s, v) - P(s', v) \mid$ | |
| Translation | $\underset{s' \in L'}{\text{argmin}} \int_v \mid P(s, v) - P(s', v) \mid$ | |
| Common sense reasoning | | |
| Planning | | |
| ... | | |

# Statistical machine translation

# Statistical machine translation

Sam was happy

# Statistical machine translation

Sam was happy

parallel corpus

# Statistical machine translation

Sam was happy

$\downarrow$ parallel corpus

Sam a fost fericit

Сэм был счастлив

# Statistical machine translation

Sam was happy

| parallel corpus

Sam a fost fericit**a**

Сэм был**а** счастлив**а**

# Statistical machine translation

Sam was happy

parallel corpus

Sam a fost fericit<span style="color:red">a</span>

Сэм был<span style="color:red">а</span> счастлив<span style="color:red">а</span>

In Thai you specify siblings by age not gender.

# Statistical machine translation

Sam was happy

$\downarrow$ parallel corpus

Sam a fost fericit<span style="color:red">a</span>

Сэм был<span style="color:red">а</span> счастлив<span style="color:red">а</span>

In Thai you specify siblings by age not gender.

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

# Statistical machine translation

Sam was happy

$\downarrow$ parallel corpus

Sam a fost fericit<span style="color:red">a</span>

Сэм был<span style="color:red">а</span> счастлив<span style="color:red">а</span>

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.

# Statistical machine translation

Sam was happy

parallel corpus

Sam a fost fericit<span style="color:red">a</span>

Сэм был<span style="color:red">а</span> счастлив<span style="color:red">а</span>

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.

# Statistical machine translation

Sam was happy

parallel corpus

Sam a fost fericita
Сэм была счастлива

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*".

# Statistical machine translation

Sam was happy

parallel corpus

Sam a fost fericit**a**

Сэм был**а** счастлив**а**

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

# Translation by imagination

French sentence

# Translation by imagination

French sentence

sample

# Translation by imagination

French sentence
$\downarrow$ sample
videos

# Translation by imagination

French sentence

$\quad\quad$ | sample
$\quad\quad$ ↓

videos

$\quad\quad$ | English generation
$\quad\quad$ ↓

# Translation by imagination

French sentence

| sample

videos

| English generation

translation

| | | |
|---|---|---|
| Recognition | $P(\text{sentence}, \text{video})$ | Narayanaswamy *et al.* 2014 |
| Retrieval | $\underset{v \in V}{\arg\max}\ P(s, v)$ | Barret *et al.* 2016 |
| Generation | $\underset{s \in L}{\arg\max}\ P(s, v)$ | Yu *et al.* 2015, Narayanaswamy *et al.* 2 |
| Question answering | $\underset{s \in L}{\arg\max}\ P(\text{Q}(s, s_q), v)$ | Barbu *et al.* in prep. |
| Disambiguation | $\underset{i \in \text{parser}(s)}{\arg\max}\ P(i, v)$ | Berzak *et al.* 2015 |
| Language acquisition | $\underset{\theta}{\arg\max} \prod_{s,v} P(s(\theta), v)$ | Yu *et al.* 2015 |
| Follow commands | $\underset{p}{\arg\max} \int_{v^+} P(C(s), v^+ v)\, E(v^+, p, v)$ | Paul *et al.* 2017 |
| Paraphrasing | $\int_v \lvert\, P(s, v) - P(s', v)\,\rvert$ | |
| Translation | $\underset{s' \in L'}{\arg\min} \int_v \lvert\, P(s, v) - P(s', v)\,\rvert$ | |
| Common sense reasoning | | |
| Planning | | |
| … | | |

| | | |
|---|---|---|
| Recognition | $P(\text{sentence}, \text{video})$ | Narayanaswamy *et al.* 2014 |
| Retrieval | $\underset{v \in V}{\text{argmax}}\ P(s, v)$ | Barret *et al.* 2016 |
| Generation | $\underset{s \in L}{\text{argmax}}\ P(s, v)$ | Yu *et al.* 2015, Narayanaswamy *et al.* 2 |
| Question answering | $\underset{s \in L}{\text{argmax}}\ P(\text{Q}(s, s_q), v)$ | Barbu *et al.* in prep. |
| Disambiguation | $\underset{i \in \text{parser}(s)}{\text{argmax}}\ P(i, v)$ | Berzak *et al.* 2015 |
| Language acquisition | $\underset{\theta}{\text{argmax}} \prod_{s,v} P(s(\theta), v)$ | Yu *et al.* 2015 |
| Follow commands | $\underset{p}{\text{argmax}} \int_{v^+} P(C(s), v^+ v)\ E(v^+, p, v)$ | Paul *et al.* 2017 |
| Paraphrasing | $\int_v \mid P(s, v) - P(s', v) \mid$ | |
| Translation | $\underset{s' \in L'}{\text{argmin}} \int_v \mid P(s, v) - P(s', v) \mid$ | |
| Common sense reasoning | $\underset{s \in L}{\text{argmax}} \int_v P(s_q, v)\ P(\text{Q}(s, s_q), v)$ | |
| Planning | | |

...

# Common sense reasoning

The trophy doesn't fit on the shelf.

# Common sense reasoning

The trophy doesn't fit on the shelf because it's too small.

Levesque *et al.* 2011

# Common sense reasoning

The trophy doesn't fit on the shelf because it's too <span style="color:red">small</span>.

Levesque *et al.* 2011

# Common sense reasoning

The trophy doesn't fit on the shelf because it's too small.

Levesque *et al.* 2011

# Common sense reasoning

The trophy doesn't fit on the shelf because it's too small.

# Common sense reasoning

The trophy doesn't fit on the shelf because it's too small.
The trophy doesn't fit on the shelf because it's too large.

Levesque *et al.* 2011

# Common sense reasoning

The trophy doesn't fit on the shelf because it's too small.
The trophy doesn't fit on the shelf because it's too large.

Levesque *et al.* 2011

# Common sense reasoning

The trophy doesn't fit on the shelf because it's too small.
The trophy doesn't fit on the shelf because it's too large.

Levesque *et al.* 2011

# Common sense reasoning

The trophy doesn't fit on the shelf because it's too small.
The trophy doesn't fit on the shelf because it's too large.

Levesque *et al.* 2011

# Common sense reasoning with vision

# Common sense reasoning with vision

$$\frac{s}{s'}$$

# Common sense reasoning with vision

$$\frac{s}{s'}$$

$s$

$s'$

# Common sense reasoning with vision

$$\frac{s}{s'}$$

YouTube

$s \longrightarrow$ 

⋮

YouTube

$s' \longrightarrow$ 

⋮

# Common sense reasoning with vision

$$\frac{s}{s'}$$



YouTube

$s \longrightarrow$ Videos $v$

YouTube

$s' \longrightarrow$ Videos $v'$

# Common sense reasoning with vision

$$\frac{s}{s'}$$



YouTube

$s \longrightarrow$ Videos $v$

YouTube

$s' \longrightarrow$ Videos $v'$

$$\frac{\displaystyle\sum_{v} S(s, v)}{\displaystyle\sum_{v'} S(s', v')}$$

# Common sense reasoning with vision

$$\frac{s}{s'}$$



$$\frac{\sum\limits_{v} S(s, v)}{\sum\limits_{v'} S(s', v')}$$

# Common sense reasoning with vision

$$\frac{s}{s'}$$

$$s \xrightarrow{\text{sample}} \text{Videos } v \qquad s' \xrightarrow{\text{sample}} \text{Videos } v'$$

$$\underbrace{\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad}$$

$$\frac{\displaystyle\int_v S(s, v)}{\displaystyle\int_{v'} S(s', v')}$$

| | | |
|---|---|---|
| Recognition | $P(\text{sentence}, \text{video})$ | Narayanaswamy *et al.* 2014 |
| Retrieval | $\underset{v \in V}{\text{argmax}}\ P(s, v)$ | Barret *et al.* 2016 |
| Generation | $\underset{s \in L}{\text{argmax}}\ P(s, v)$ | Yu *et al.* 2015, Narayanaswamy *et al.* 2 |
| Question answering | $\underset{s \in L}{\text{argmax}}\ P(\text{Q}(s, s_q), v)$ | Barbu *et al.* in prep. |
| Disambiguation | $\underset{i \in \text{parser}(s)}{\text{argmax}}\ P(i, v)$ | Berzak *et al.* 2015 |
| Language acquisition | $\underset{\theta}{\text{argmax}} \prod_{s,v} P(s(\theta), v)$ | Yu *et al.* 2015 |
| Follow commands | $\underset{p}{\text{argmax}} \int_{v^+} P(\text{C}(s), v^+ v)\ E(v^+, p, v)$ | Paul *et al.* 2017 |
| Paraphrasing | $\int_v \mid P(s, v) - P(s', v) \mid$ | |
| Translation | $\underset{s' \in L'}{\text{argmin}} \int_v \mid P(s, v) - P(s', v) \mid$ | |
| Common sense reasoning | $\underset{s \in L}{\text{argmax}} \int_v P(s_q, v)\ P(\text{Q}(s, s_q), v)$ | |
| Planning | | |

…

| | | |
|---|---|---|
| Recognition | $P(\text{sentence}, \text{video})$ | Narayanaswamy *et al.* 2014 |
| Retrieval | $\underset{v \in V}{\text{argmax}}\ P(s, v)$ | Barret *et al.* 2016 |
| Generation | $\underset{s \in L}{\text{argmax}}\ P(s, v)$ | Yu *et al.* 2015, Narayanaswamy *et al.* 2 |
| Question answering | $\underset{s \in L}{\text{argmax}}\ P(\text{Q}(s, s_q), v)$ | Barbu *et al.* in prep. |
| Disambiguation | $\underset{i \in \text{parser}(s)}{\text{argmax}}\ P(i, v)$ | Berzak *et al.* 2015 |
| Language acquisition | $\underset{\theta}{\text{argmax}} \prod_{s,v} P(s(\theta), v)$ | Yu *et al.* 2015 |
| Follow commands | $\underset{p}{\text{argmax}} \int_{v^+} P(C(s), v^+ v)\, E(v^+, p, v)$ | Paul *et al.* 2017 |
| Paraphrasing | $\int_v \| P(s, v) - P(s', v) \|$ | |
| Translation | $\underset{s' \in L'}{\text{argmin}} \int_v \| P(s, v) - P(s', v) \|$ | |
| Common sense reasoning | $\underset{s \in L}{\text{argmax}} \int_v P(s_q, v)\, P(\text{Q}(s, s_q), v)$ | |
| Planning | $\underset{s \in L}{\text{argmax}} \int_v P(s, v_0 : v : v_n)$ | |
| … | | |

# 5g of brain

# 5g of brain
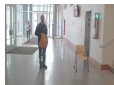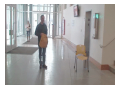
# Planning

# Planning

# Planning



Danny carried the backpack to the chair.

# Planning



Danny carried the backpack to the chair.

# Planning



Danny carried the backpack to the chair.

# Planning



Danny carried the backpack to the chair.

# Planning



Danny carried the backpack to the chair.

| | | |
|---|---|---|
| Recognition | $P(\text{sentence}, \text{video})$ | Narayanaswamy *et al.* 2014 |
| Retrieval | $\underset{v \in V}{\text{argmax}}\; P(s, v)$ | Barret *et al.* 2016 |
| Generation | $\underset{s \in L}{\text{argmax}}\; P(s, v)$ | Yu *et al.* 2015, Narayanaswamy *et al.* 2 |
| Question answering | $\underset{s \in L}{\text{argmax}}\; P(\text{Q}(s, s_q), v)$ | Barbu *et al.* in prep. |
| Disambiguation | $\underset{i \in \text{parser}(s)}{\text{argmax}}\; P(i, v)$ | Berzak *et al.* 2015 |
| Language acquisition | $\underset{\theta}{\text{argmax}} \prod_{s,v} P(s(\theta), v)$ | Yu *et al.* 2015 |
| Follow commands | $\underset{p}{\text{argmax}} \int_{v^+} P(C(s), v^+ v)\, E(v^+, p, v)$ | Paul *et al.* 2017 |
| Paraphrasing | $\int_v \mid P(s, v) - P(s', v) \mid$ | |
| Translation | $\underset{s' \in L'}{\text{argmin}} \int_v \mid P(s, v) - P(s', v) \mid$ | |
| Common sense reasoning | $\underset{s \in L}{\text{argmax}} \int_v P(s_q, v)\; P(\text{Q}(s, s_q), v)$ | |
| Planning | $\underset{s \in L}{\text{argmax}} \int_v P(s, v_0 : v : v_n)$ | |
| ... | | |

# 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

# A person standing in front of a stove

# 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

# The long road ahead . . .

Coherent stories

3D

Physics: Forces & contact relations

Segmentation

Parts and low-level features

Theory of mind

Social understanding

# Social understanding

# Social understanding

# 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

# The long road ahead . . .

Coherent stories

3D

Physics: Forces & contact relations

Segmentation

Parts and low-level features

Theory of mind

Social understanding

Modification

# The long road ahead . . .

Coherent stories

3D

Physics: Forces & contact relations

Segmentation

Parts and low-level features

Theory of mind

Social understanding

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.

## The long road ahead . . .

Coherent stories

3D

Physics: Forces & contact relations

Segmentation

Parts and low-level features

Theory of mind

Social understanding

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

# The long road ahead . . .

Coherent stories

3D

Physics: Forces & contact relations

Segmentation

Parts and low-level features

Theory of mind

Social understanding

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, Candace Ross, Yen-Ling Kuo, Jonathan Malmaud
Daniel Harari, Battushig Myanganbayar, David Mayo, Nazar Ilamanov
Boris Katz, Shimon Ullman, Josh Tenenbaum

Siddharth Narayanaswamy, Jeffrey Siskind
Victor Carbarera, Santiago Perez
Sven Dickinson, Song Wang
Daniel Barrett, Haonan Yu
Maria Ryskina, Sergey Voronov

| | | |
|---|---|---|
| Recognition | $P(\text{sentence}, \text{video})$ | Narayanaswamy *et al.* 2014 |
| Retrieval | $\underset{v \in V}{\text{argmax}}\ P(s, v)$ | Barret *et al.* 2016 |
| Generation | $\underset{s \in L}{\text{argmax}}\ P(s, v)$ | Yu *et al.* 2015, Narayanaswamy *et al.* 2 |
| Question answering | $\underset{s \in L}{\text{argmax}}\ P(\text{Q}(s, s_q), v)$ | Barbu *et al.* in prep. |
| Disambiguation | $\underset{i \in \text{parser}(s)}{\text{argmax}}\ P(i, v)$ | Berzak *et al.* 2015 |
| Language acquisition | $\underset{\theta}{\text{argmax}} \prod_{s,v} P(s(\theta), v)$ | Yu *et al.* 2015 |
| Follow commands | $\underset{p}{\text{argmax}} \int_{v^+} P(C(s), v^+ v)\, E(v^+, p, v)$ | Paul *et al.* 2017 |
| Paraphrasing | $\int_{v} \mid P(s, v) - P(s', v) \mid$ | |
| Translation | $\underset{s' \in L'}{\text{argmin}} \int_{v} \mid P(s, v) - P(s', v) \mid$ | |
| Common sense reasoning | $\underset{s \in L}{\text{argmax}} \int_{v} P(s_q, v)\, P(\text{Q}(s, s_q), v)$ | |
| Planning | $\underset{s \in L}{\text{argmax}} \int_{v} P(s, v_0 : v : v_n)$ | |
| … | | |