## Structure and Interpretation of Neural Codes







Jacob Andreas

## Translating Neuralese





#### Jacob Andreas, Anca Drăgan and Dan Klein





[Wagner et al. 03, Sukhbaatar et al. 16, Foerster et al. 16]

#### Learning to Communicate







#### Learning to Communicate

4







#### Neuralese

1.02.3-0.30.4-1.21.1 X







#### Translating neuralese













#### Interoperate with autonomous systems

#### • **Diagnose** errors

#### • Learn from solutions

[Lazaridou et al. 16]

#### Translating neuralese







# Natural language & neuralese **Statistical** machine translation **Semantic** machine translation Implementation details Evaluation







## Natural language & neuralese Statistical machine translation

#### **Semantic** machine translation

Implementation details Evaluation







# Natural language & neuralese **Statistical** machine translation **Semantic** machine translation Implementation details Evaluation







# Natural language & neuralese Statistical machine translation Semantic machine translation

## Implementation details

Evaluation



# Natural language & neuralese Statistical machine translation **Semantic** machine translation Implementation details

Evaluation





# Natural language & neuralese **Statistical** machine translation

#### **Semantic** machine translation

# Implementation details

Evaluation









[e.g. Koehn 10]

#### A statistical MT problem

all clear











#### How do we induce a translation model?

#### A statistical MT problem









## $\max p([o] | a]) p([a])$ а $\propto \max \sum p([o]) + p([a]) + p([a]) + p([a])$

#### A statistical MT problem









































## Σp(, I not sure) p( not sure)







#### Stat MT criterion doesn't capture meaning









## Natural language & neuralese **X** Statistical machine translation

#### **Semantic** machine translation

Implementation details Evaluation





#### The meaning of an utterance is given by its truth conditions





[Davidson 67]











#### The meaning of an utterance is given by its truth conditions







[Davidson 67]















#### The meaning of an utterance is given by its truth conditions







#### (loc (goal blue) north)







#### The meaning of an utterance is given by its truth conditions the distribution over states in which it is uttered







[Beltagy et al. 14]



0.2



0.001









#### The meaning of an utterance is given by its truth conditions

#### the distribution over states in which it is uttered

#### the **belief** it induces in listeners





0.4



## A "semantic MT" problem



0.2











# The meaning of an utterance is given by

- the distribution over states in which it is uttered
  - or equivalently, the **belief** it induces in listeners







# The meaning of an utterance is given by

vector rather than a sequence of tokens.

- the distribution over states in which it is uttered
  - or equivalently, the **belief** it induces in listeners

This distribution is well-defined even if the "utterance" is a



















# In the intersection






























### Interlingua!

















# argmin $KL(\beta(0)) || \beta(a))$

















## argmin $KL(\beta(0) || \beta(a))$

### Computing representations





## argmin $\mathbb{A}$ KL( $\beta(\mathbb{O}) \parallel \beta(\mathbb{A})$ )



### Computing representations: sparsity



## argmin $KL(\beta(0) || \beta(a))$

agent policy









44

### agent model





## argmin $\mathcal{KL}(\beta(0) | \beta(a))$



### Computing representations: smoothing





## argmin $KL(\beta(0) || \beta(a))$

human









## argmin $KL(\beta(0) || \beta(a))$

### human policy



human model



## argmin $\mathbb{A}$ KL( $\beta(\mathbb{O}) \parallel \beta(\mathbb{O})$ )







## argmin $\mathbb{A}$ KL( $\beta(\mathbb{P}) \parallel \beta(\mathbb{P})$ )

### Computing KL







## argmin $KL(\beta(0) | \beta(a))$

### Computing KL









## $\operatorname{argmin}_{a} \operatorname{KL}(\beta(0) | \beta(a))$

### Computing KL: sampling







## argmin $KL(\beta(0) || \beta(a))$

### Finding translations



### Finding translations: brute force



## argmin $KL(\beta(0) || \beta(a))$

after you





### Finding translations: brute force

## argmin $KL(\beta(0) || \beta(a))$

### going north —

l'm done

atter you







### Finding translations







## Natural language & neuralese Statistical machine translation **Semantic** machine translation Implementation details

Evaluation

### Outline



### Referring expression games























### Evaluation: translator-in-the-loop











### Evaluation: translator-in-the-loop





























### English → English\*











































### magenta, hot, violet









### magenta, hot, violet

olive, puke, pea







### Ø magenta, hot, rose

### magenta, hot, violet

### olive, puke, pea

pinkish, grey, dull









### Experiment: image references





small brown, light brown, dark brown

large bird, black wings, black crown











### Statistical MT



### Semantic MT



### Experiment: driving game

### Neuralese $\rightarrow$ Neuralese

Neuralese ↔ English\*




#### How to translate



#### at goal done left to top





you first following going down

going in intersection proceed going



#### Classical notions of "meaning" apply even to un-language-like things (e.g. RNN states)

• These meanings can be compactly represented without logical forms if we have access to world states

Communicating policies "say" interpretable things!





#### Classical notions of "meaning" apply even to non-language-like things (e.g. RNN states)

# These meanings can be compactly represented

without logical forms if we have access to world states

Communicating policies "say" interpretable things!





 Classical notions of "meaning" apply even to non-language-like things (e.g. RNN states)

 These meanings can be compactly represented without logical forms if we have access to world states

Communicating policies "say" interpretable things!





# $\operatorname{argmin}_{a} \operatorname{KL}(\beta(\mathfrak{G}) | \beta(\mathfrak{G}))$

#### Limitations

 $KL(p \parallel q) = \sum_{i} p(\sum_{i}) \log \frac{p(\sum_{i})}{q(\sum_{i})}$ 



### but what about compositionality?

## Analogs of linguistic structure in deep representations







#### Jacob Andreas and Dan Klein





#### at goal done



#### "Flat" semantics



you first following

going in intersection proceed going





































[FitzGerald et al. 2013]



everything but the blue shapes orange square and non-squares











#### [FitzGerald et al. 2013]



#### lambda x: not(blue(x)) lambda x: or(orange(x), not(square(x))











































1.0 -0.3







































































#### Translation criterion



## $q(\ensuremath{\textcircled{o}}\ensuremath{,}\ensuremath{\textcircled{o}}\ensuremath{,}\ensuremath{\textcircled{o}}\ensuremath{,}\ensuremath{\textcircled{o}}\ensuremath{,}\ensuremath{\textcircled{o}}\ensuremath{,}\ensuremath{\textcircled{o}}\ensuremath{,}\ensuremath{\textcircled{o}}\ensuremath{,}\ensuremath{\textcircled{o}}\ensuremath{,}\ensuremath{\textcircled{o}}\ensuremath{,}\ensuremath{\textcircled{o}}\ensuremath{,}\ensuremath{\textcircled{o}}\ensuremath{,}\ensuremath{\textcircled{o}}\ensuremath{,}\ensuremath{\textcircled{o}}\ensuremath{,}\ensuremath{\textcircled{o}}\ensuremath{,}\ensuremath{\textcircled{o}}\ensuremath{,}\ensuremath{\textcircled{o}}\ensuremath{,}\ensuremath{e}\ensuremath{,}\ensuremath{e}\ensuremath{,}\ensuremath{e}\ensuremath{,}\ensuremath{e}\ensuremath{e}\ensuremath{,}\ensuremath{e}\ensuremath{,}\ensuremath{e}\ensuremath{,}\ensuremath{e}\ensuremath{,}\ensuremath{e}\ensuremath{,}\ensuremath{e}\ensuremath{,}\ensuremath{e}\ensuremath{,}\ensuremath{e}\ensuremath{,}\ensuremath{e}\ensuremath{,}\ensuremath{e}\ensuremath{,}\ensuremath{e}\ensuremath{,}\ensuremath{e}\ensuremath{,}\ensuremath{e}\ensuremath{,}\ensuremath{e}\ensuremath{,}\ensuremath{,}\ensuremath{e}\ensuremath{e}\ensuremath{,}\ensuremath{e}\ensuremath{,}\ensuremath{e}\ensuremath{,}\ensuremath{e}\ensuremath{,}\ensuremath{e}\ensuremath{,}\ensuremath{e}\ensuremath{,}\ensuremath{e}\ensuremath{,}\ensuremath{e}\ensuremath{,}\ensuremath{e}\ensuremath{,}\ensuremath{e}\ensuremath{,}\ensuremath{e}\ensuremath{,}\ensuremath{e}\ensuremath{,}\ensuremath{e}\ensuremath{,}\ensuremath{e}\ensuremath{,}\ensuremath{e}\ensuremath{,}\ensuremath{e}\ensuremath{,}\ensuremath{,}\ensur$





#### Translation criterion



## $q(\bigcirc, \bigcirc) = \mathbf{E}[\beta(\bigcirc) = \beta(\bigcirc)]$









# "High-level" communicative behavior "Low-level" message structure



101





## "High-level" communicative behavior "Low-level" message structure











103







# everything but squares













# everything but squares















# everything but squares







### Theories of model behavior: random

$$\begin{array}{cccc} -0.1 & 1.3 \\ 0.5 & -0.4 \\ 0.2 & 1.0 \end{array}$$









#### Theories of model behavior: literal

$$\begin{array}{cccc} -0.1 & 1.3 \\ 0.5 & -0.4 \\ 0.2 & 1.0 \end{array}$$














#### Evaluation: high-level scene agreement









#### Evaluation: high-level object agreement









# "High-level" communicative behavior "Low-level" message structure

111



## Collecting translation data

all the red shapes

blue objects

everything but red

green squares

not green squares







## Collecting translation data

$$\lambda x.red(x)$$

$$\lambda x.blu(x)$$

$$\lambda x.\neg red(x)$$





### Collecting translation data

$$\lambda x.red(x)$$

$$\lambda x.blu(x)$$

$$\lambda x.\neg red(x)$$

$$\begin{array}{c} \bullet & 0.1 & -0.3 & 0.5 & 1.3 \\ \bullet & -0.3 & 0.2 & 0.1 & 0.3 \\ \bullet & 1.4 & -0.3 & -0.5 & 0.3 \\ \bullet & 0.2 & -0.2 & 0.5 & -0.3 \\ \bullet & 0.3 & -1.3 & -1.5 & 0.3 \end{array}$$









#### Extracting related pairs











#### Extracting related pairs









## argmin

#### Learning compositional operators







## Evaluating learned operators





## Evaluating learned operators









## Evaluating learned operators

















#### Evaluation: scene agreement for negation



0

0



Input

#### Predicted

True

all the toys that of are not red all items that are only the blue and not blue or green green objects

#### Visualizing negation

every thing that is red











#### Evaluation: scene agreement for disjunction



123







## Visualizing disjunction







- We can translate between neuralese and natural lang. by grounding in distributions over world states
- Under the right conditions, neuralese exhibits interpretable pragmatics & compositional structure
- Not just communication games—language might be a good general-purpose tool for interpreting deep reprs.





- by grounding in distributions over world states
- Under the right conditions, neuralese exhibits

We can translate between neuralese and natural lang.

interpretable pragmatics & compositional structure

 Not just communication games—language might be a good general-purpose tool for interpreting deep reprs.





- We can translate between neuralese and natural lang. by grounding in distributions over world states
- Under the right conditions, neuralese exhibits interpretable pragmatics & compositional structure
- Not just communication games—language might be a good general-purpose tool for interpreting deep reprs.

#### Conclusions







#### Conclusions







http://github.com/jacobandreas/{neuralese,rnn-syn}

