

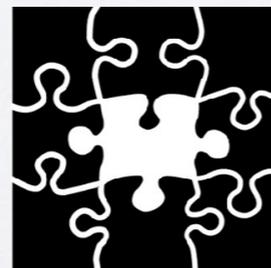
LTG Seminar | 14 March 2022

# Efficient Strategies of Language Production: An Information-Theoretic Analysis

Mario Giulianelli

Institute for Logic, Language and Computation  
University of Amsterdam

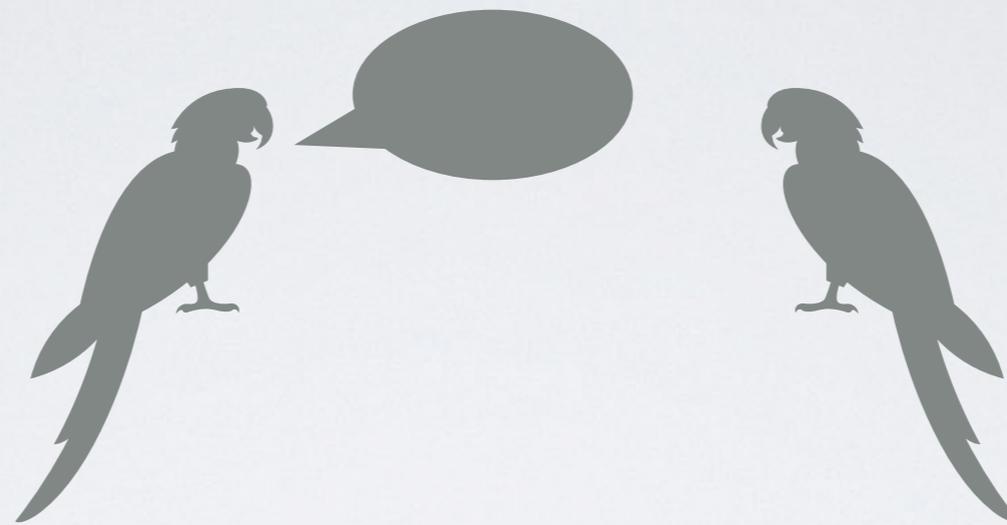
`m.giulianelli@uva.nl` | `glnmario.github.io`



# Efficient Strategies of Language Production: An Information-Theoretic Analysis

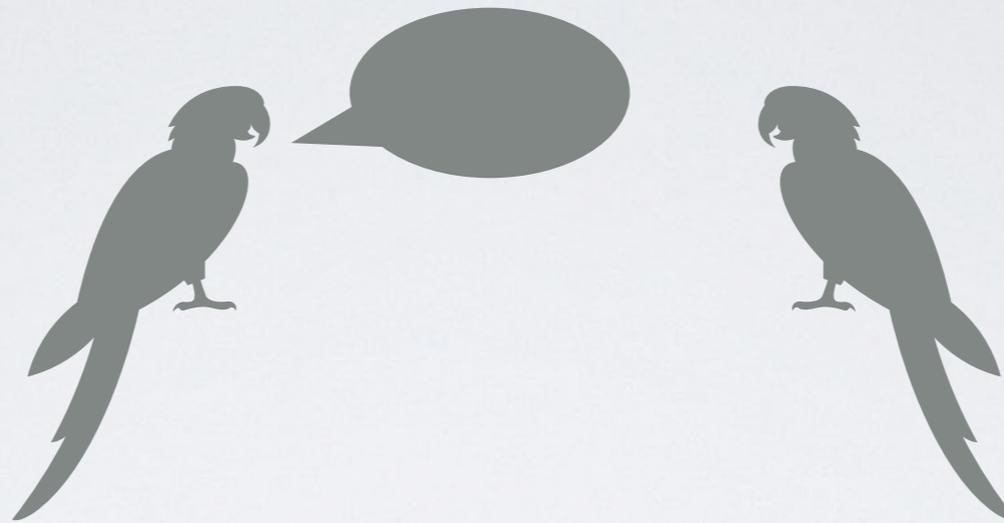
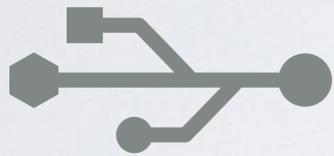
- ◆ Three studies, with Raquel Fernández and Arabella Sinclair
- ◆ An information-theoretic model of language production
- ◆ A hypothesis: humans follow efficient strategies of information transmission, leading to near-optimal collaborative effort
- ◆ Utterance surprisal (information density) as a proxy of effort, estimated with an autoregressive transformer language model

# Collaborative effort in language production



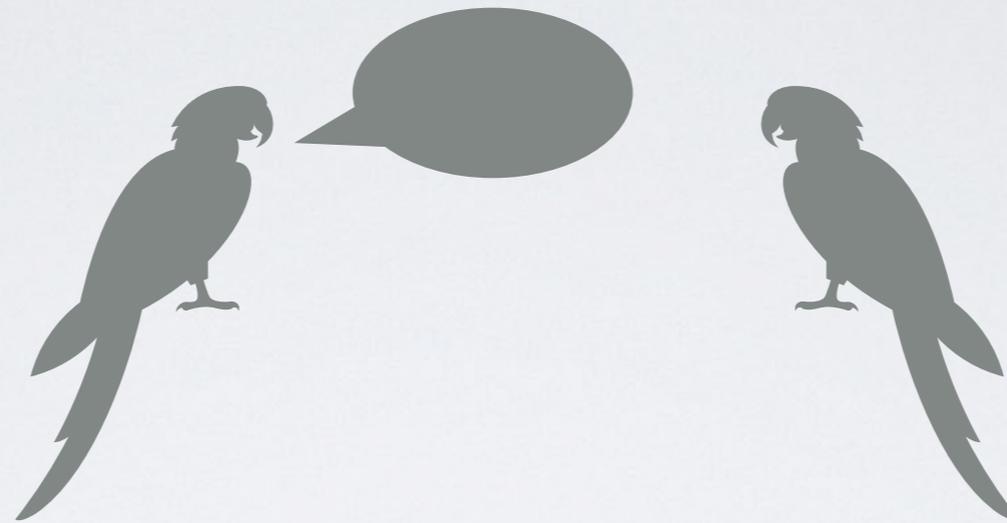
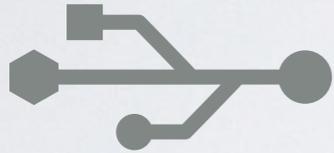
# Collaborative effort in language production

PRODUCTION  
EFFORT

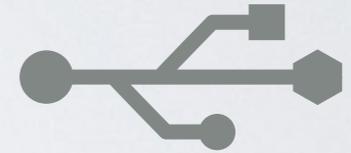


# Collaborative effort in language production

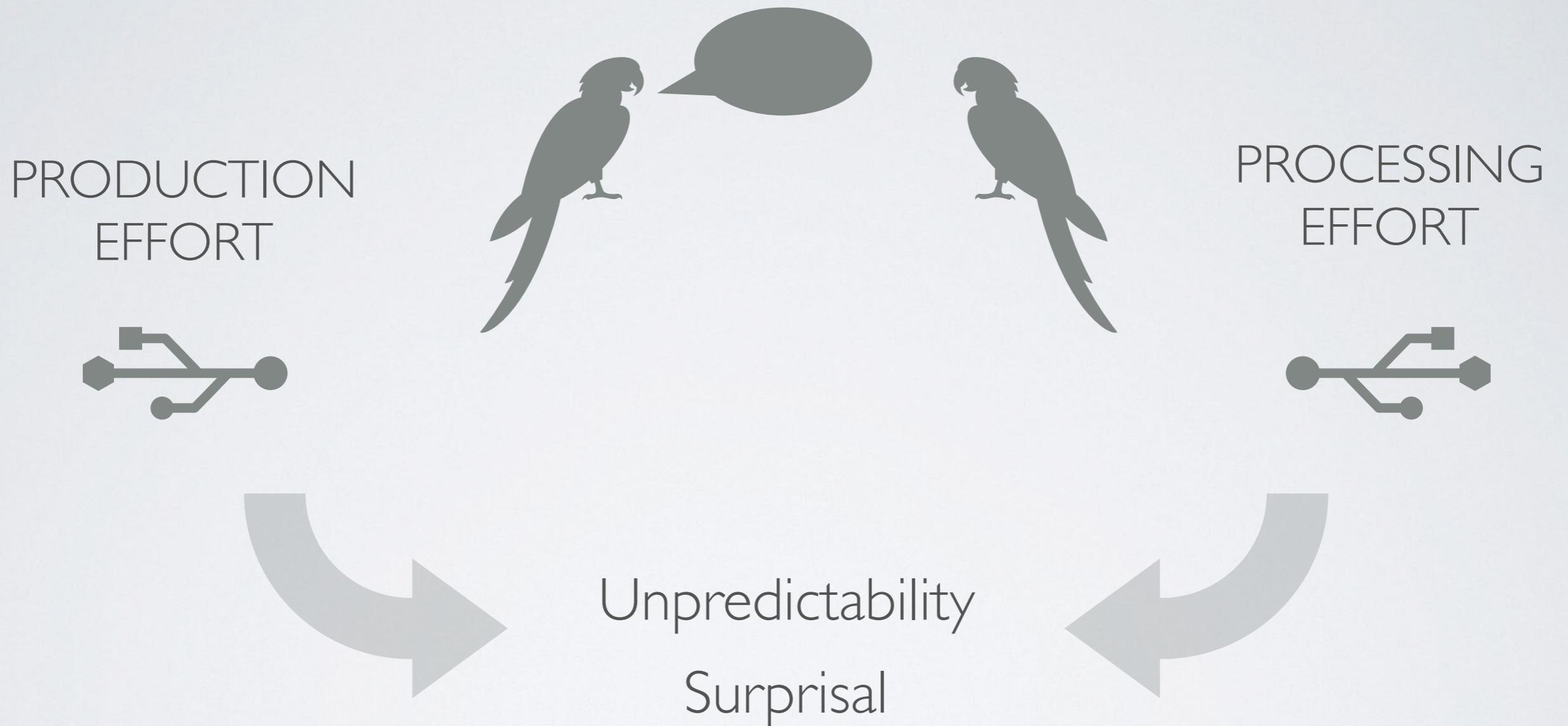
PRODUCTION  
EFFORT



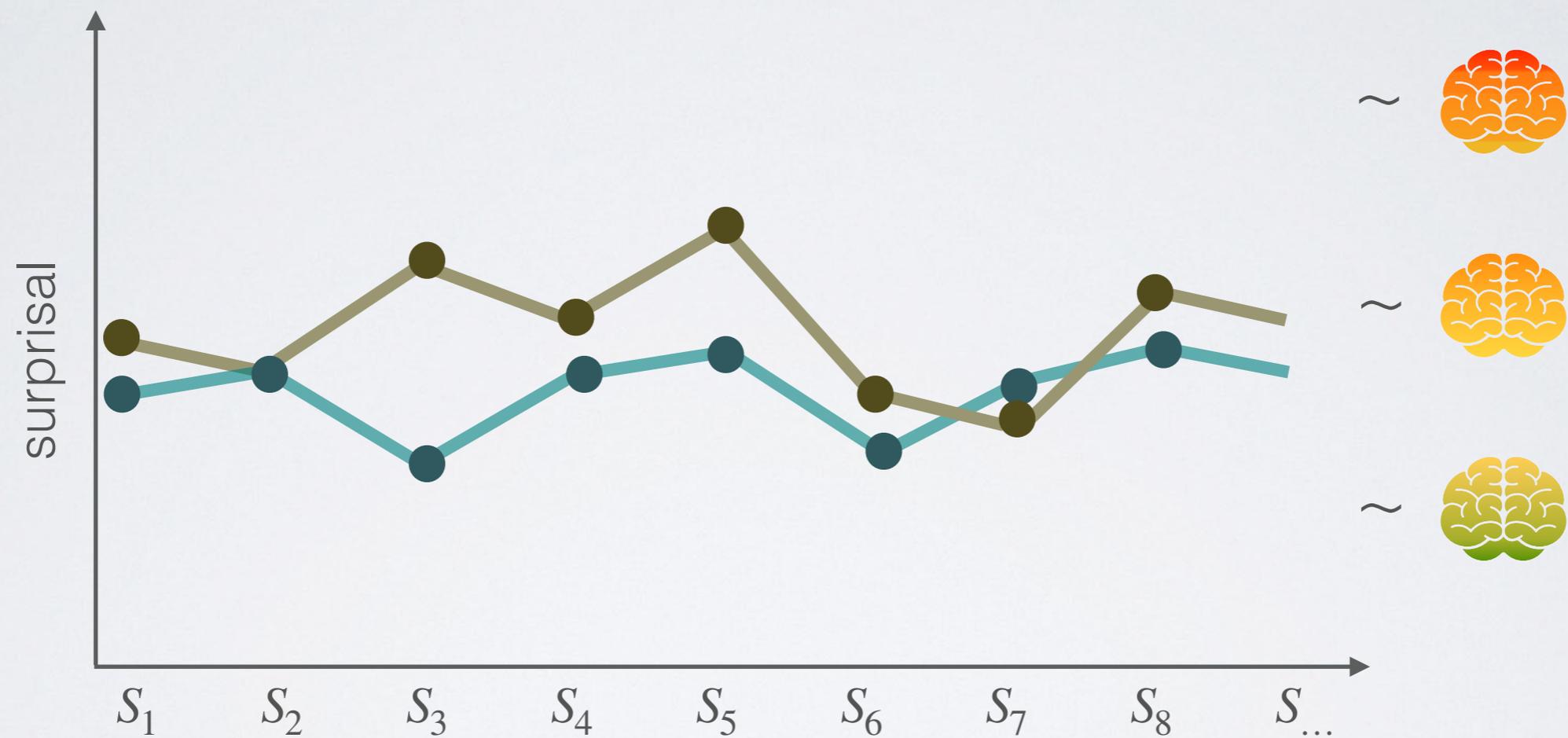
PROCESSING  
EFFORT



# Collaborative effort in language production

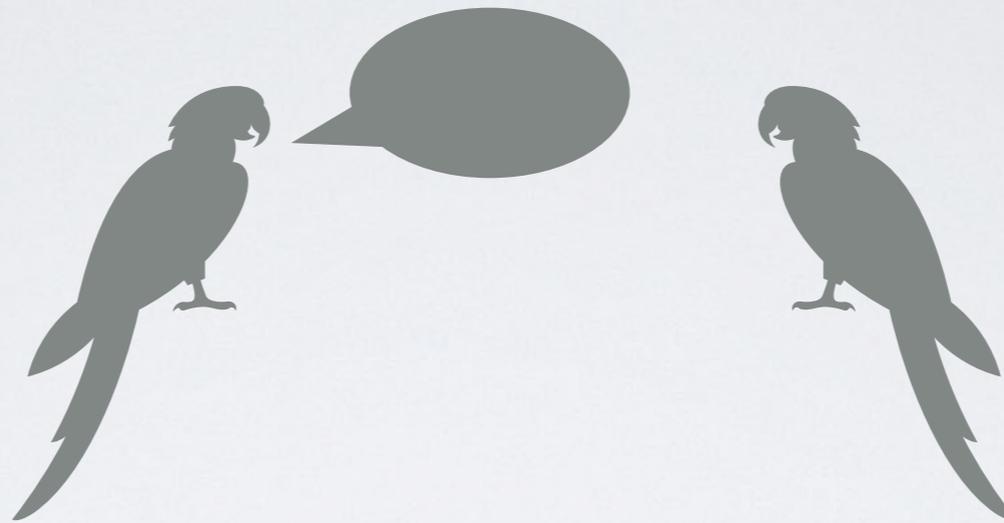
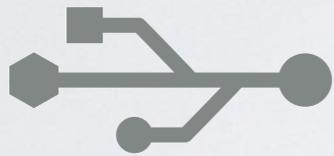


# Collaborative effort in language production

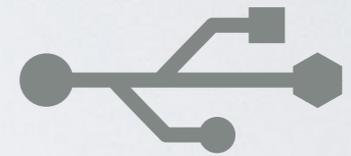


# Collaborative effort in language production

PRODUCTION  
EFFORT



PROCESSING  
EFFORT

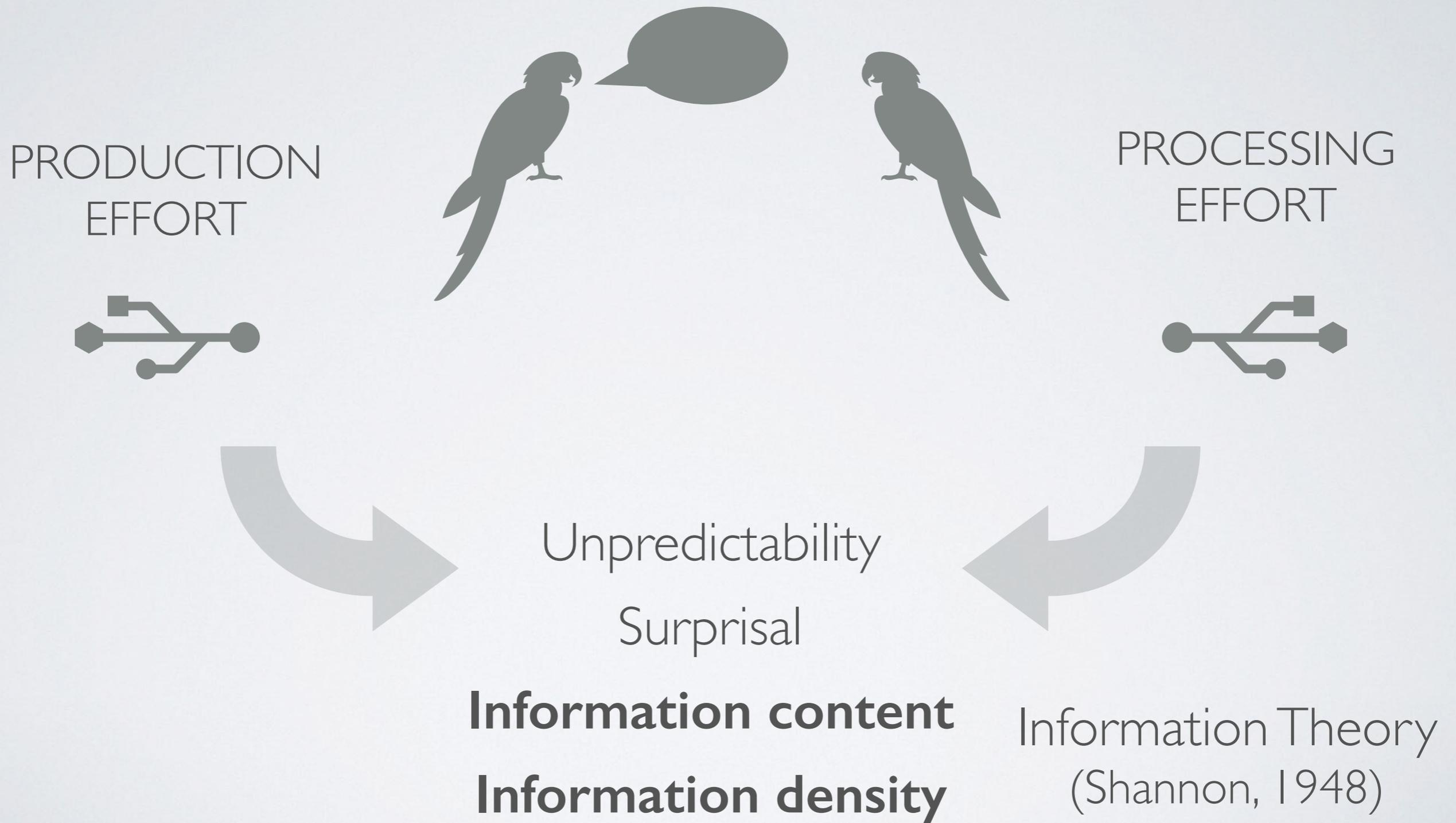


phoneme duration,  
syntactic structures,  
turn taking

Unpredictability  
Surprisal

perception,  
reading,  
sentence interpretation

# Collaborative effort in language production



# Information-theoretic measures

# Information-theoretic measures

$H(S)$

Out-of-context surprisal

$$-\log_2 P(S)$$

Amcore Financial Inc. said it agreed to acquire Central of Illinois Inc. in a stock swap.

Shareholders of Central, a bank holding company based in Sterling, will receive Amcore stock equal to 10 times Central's 1989 earnings, Amcore said.

**For the first nine months of 1989, Central earned \$2 million.**

# Information-theoretic measures

$H(S | C)$  In-context surprisal

$$-\log_2 P(S | C)$$

Amcore Financial Inc. said it agreed to acquire Central of Illinois Inc. in a stock swap.

Shareholders of Central, a bank holding company based in Sterling, will receive Amcore stock equal to 10 times Central's 1989 earnings, Amcore said.

**For the first nine months of 1989, Central earned \$2 million.**

# Information-theoretic measures

$I(S; C)$

Context informativeness

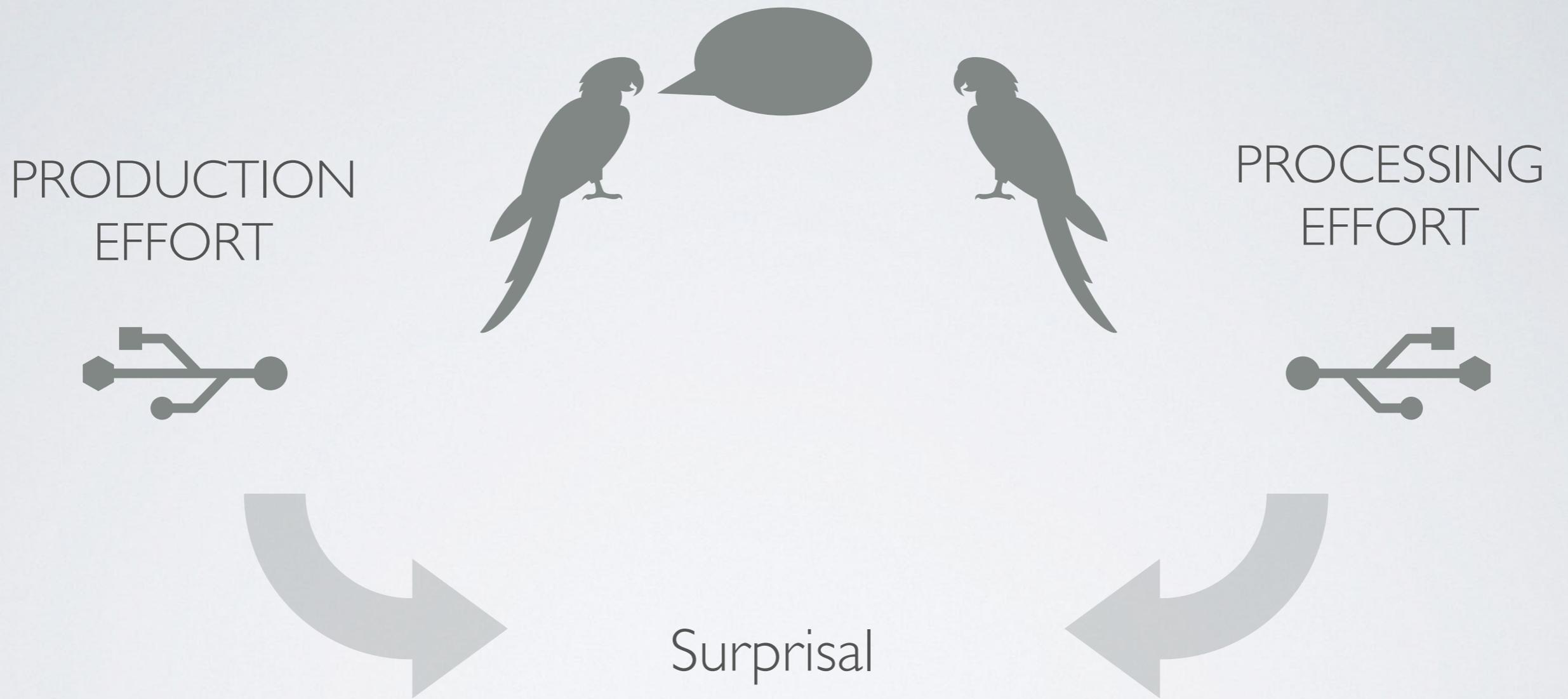
$$H(S) - H(S | C) = -\log_2 P(S) + \log_2 P(S | C)$$

Amcore Financial Inc. said it agreed to acquire Central of Illinois Inc. in a stock swap.

Shareholders of Central, a bank holding company based in Sterling, will receive Amcore stock equal to 10 times Central's 1989 earnings, Amcore said.

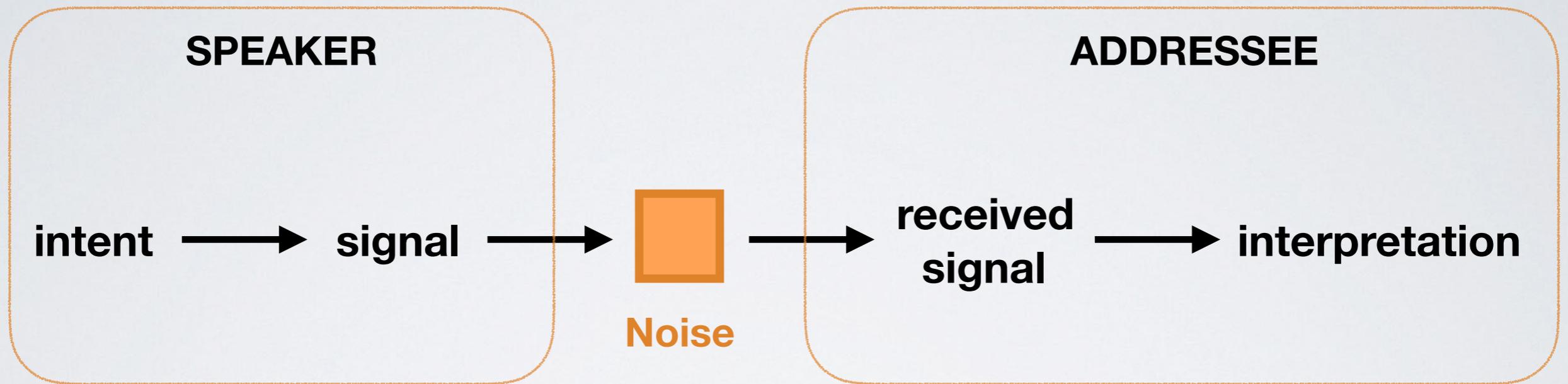
For the first nine months of 1989, Central earned \$2 million.

# Collaborative effort in language production



# **Strategies of language production**

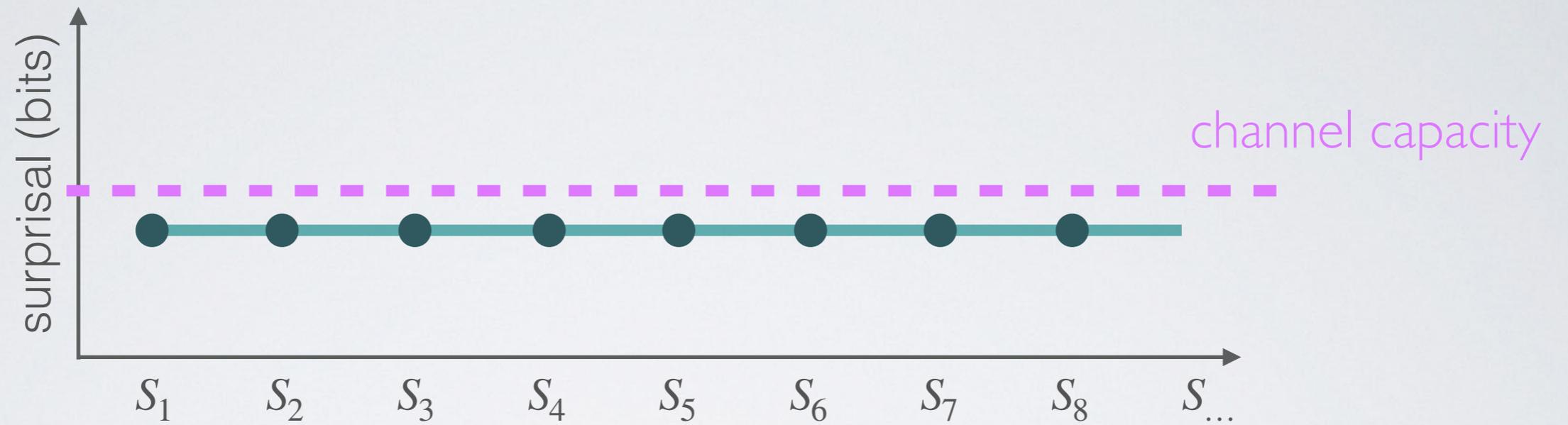
# Strategies of language production



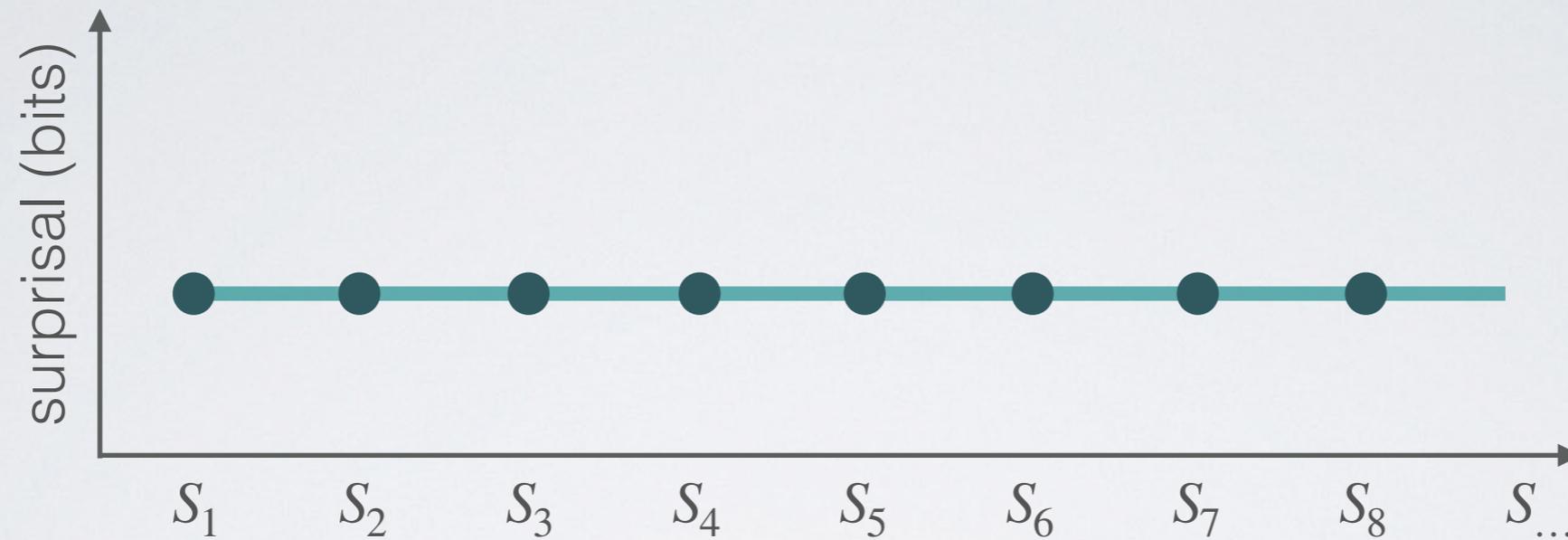
For any given degree of noise, it is possible to communicate nearly error-free up to a *maximum rate* through the channel.

(Shannon, 1948)

# Strategies of language production



# Strategies of language production

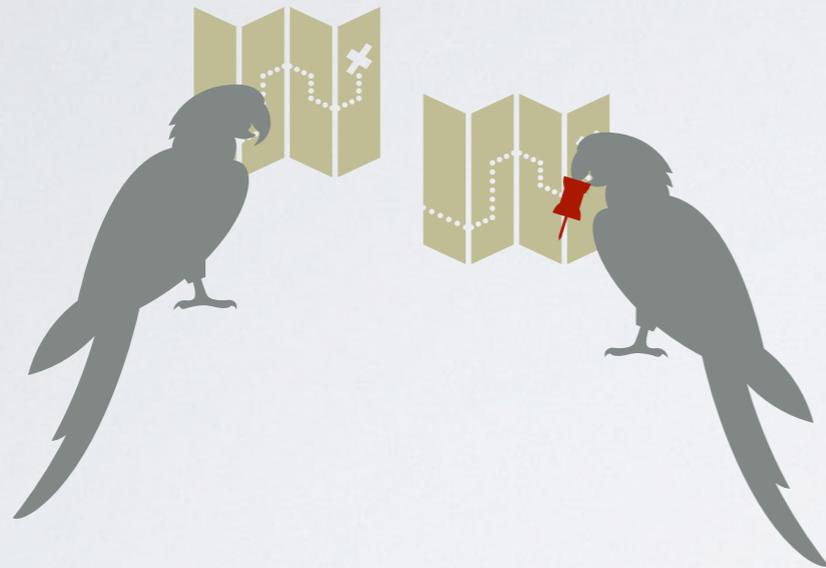


Constancy Rate Principle  
(Genzel & Charniak, 2002)

## Study I

**Does the Constancy Rate Principle hold  
in task-oriented dialogues?**

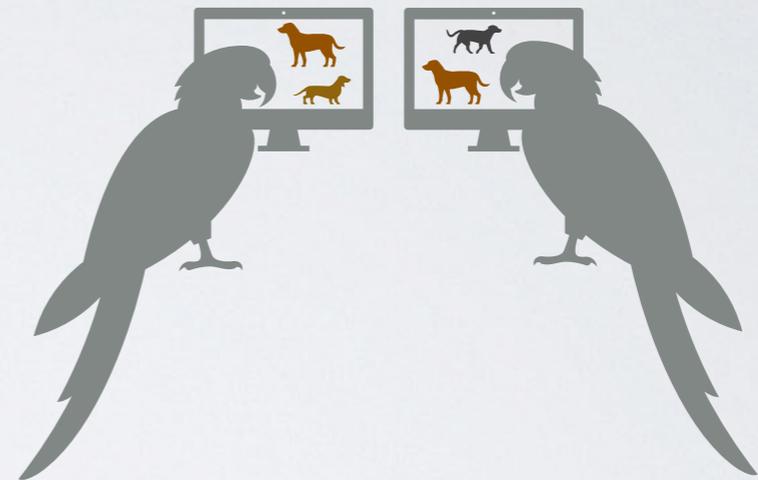
# Task-oriented dialogue corpora



## MapTask

instruction giving  
and following

(Anderson et al., 1991)

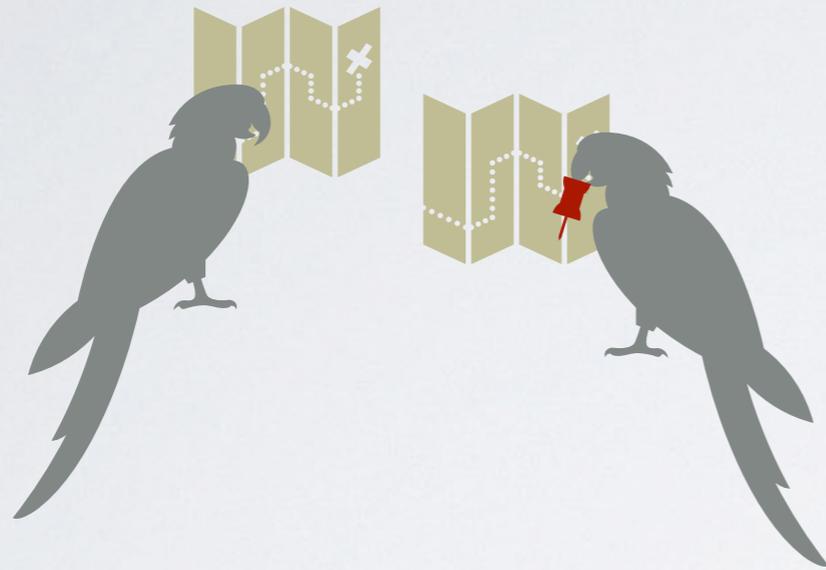


## PhotoBook

written cooperative  
reference game

(Haber et al., 2019)

# Task-oriented dialogue corpora



**MapTask**  
instruction giving  
and following  
(Anderson et al., 1991)

G	okay
G	go down
F	down towards missionary camp
G	aye
G	just right down uh-huh
F	right
F	i'm i've stopped right above it
G	then
G	you've what
F	i'm right above the missionary camp

Transaction 2

G	just go right down to round down to [...]
F	right
G	okay
F	so go round it
G	uh-huh
G	round to the left uh the left yeah
F	keep going down
G	uh-huh
G	just stop where you when you reach [...]

Transaction 3

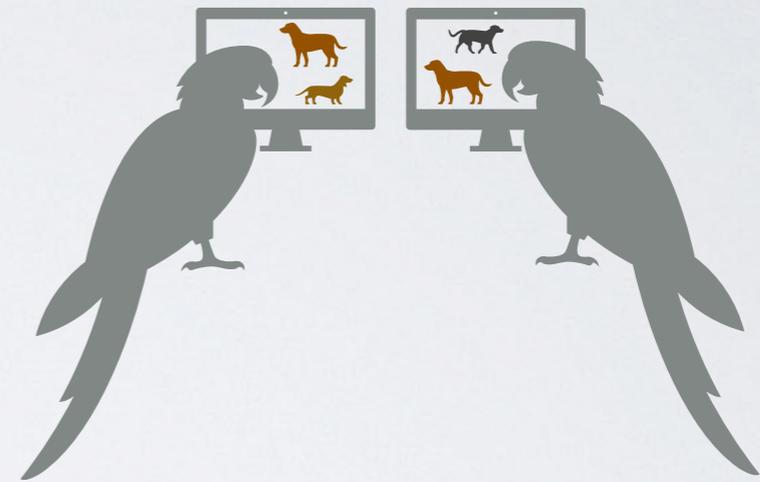
# Task-oriented dialogue corpora

## Round 3

B	pink sweater woman and man w/umbrella
A	yes
B	statue man with umbrella
A	y
B	guy in black suit with 2 plaid blue umbrellas
A	no

## Round 4

B	guy in black suit with 2 plaid blue umbrellas
A	no
B	lady in pink and guy with umbrella
A	no
B	statue/umbralla again
A	no



**PhotoBook**  
written cooperative  
reference game  
(Haber et al., 2019)

# Task-oriented dialogue corpora

Round 3

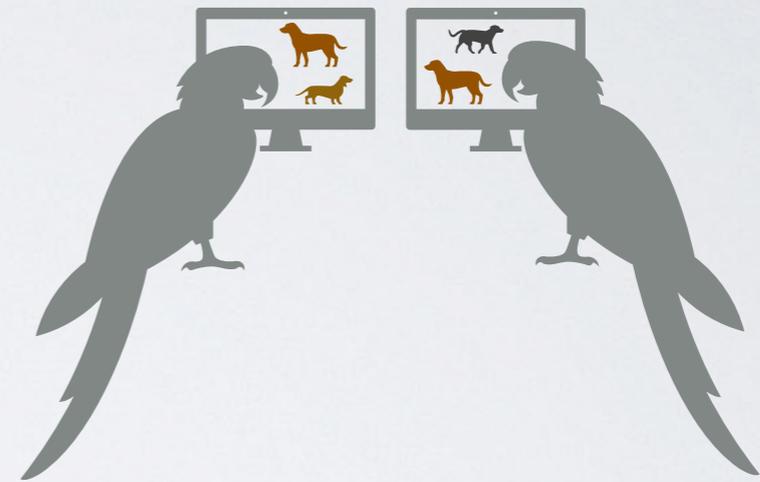
B	pink sweater woman and man w/umbrella
A	yes
B	statue man with umbrella
A	y
B	guy in black suit with 2 plaid blue umbrellas
A	no

Round 4

B	guy in black suit with 2 plaid blue umbrellas
A	no
B	lady in pink and guy with umbrella
A	no
B	statue/umbralla again
A	no

## Reference chains

A	man eating slice of pizza
B	last one for me is guy with pizza
A	pizza eater
B	pizza

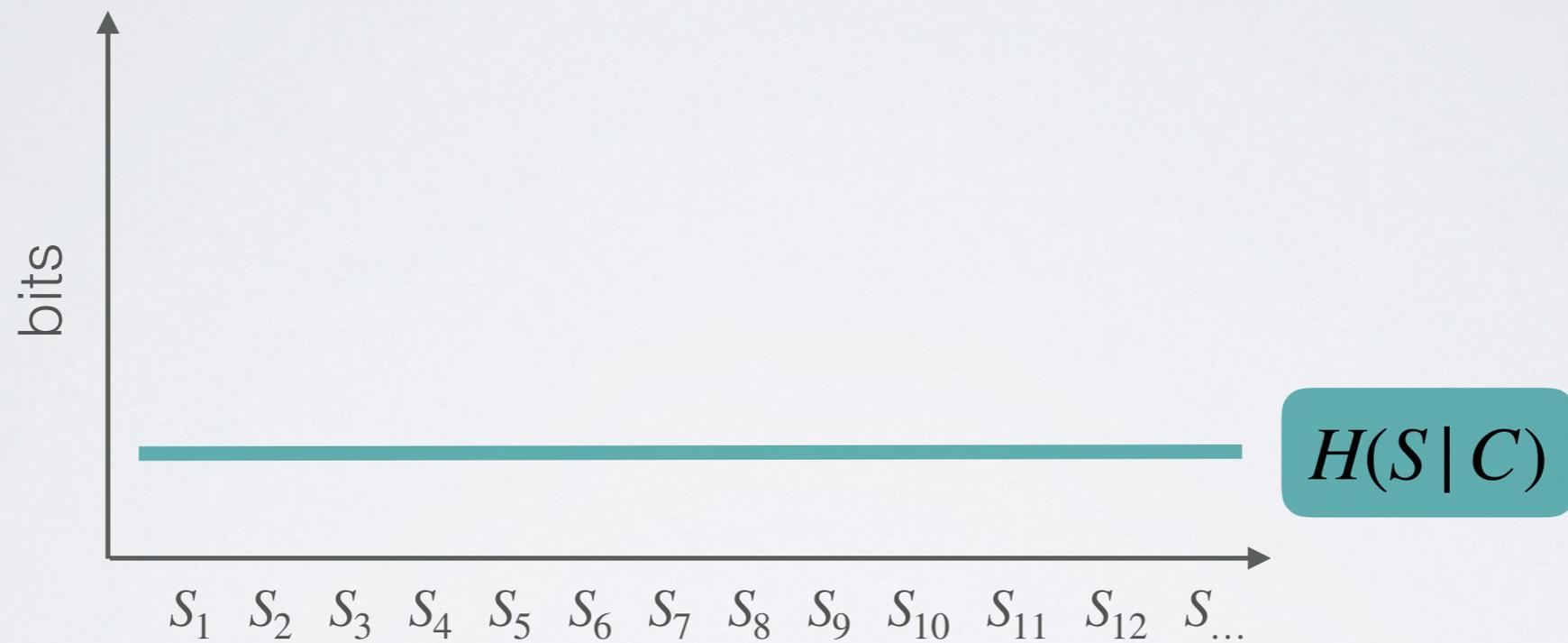


**PhotoBook**  
written cooperative  
reference game  
(Haber et al., 2019)

**Reference chain dataset**  
(Takmaz et al., 2020)

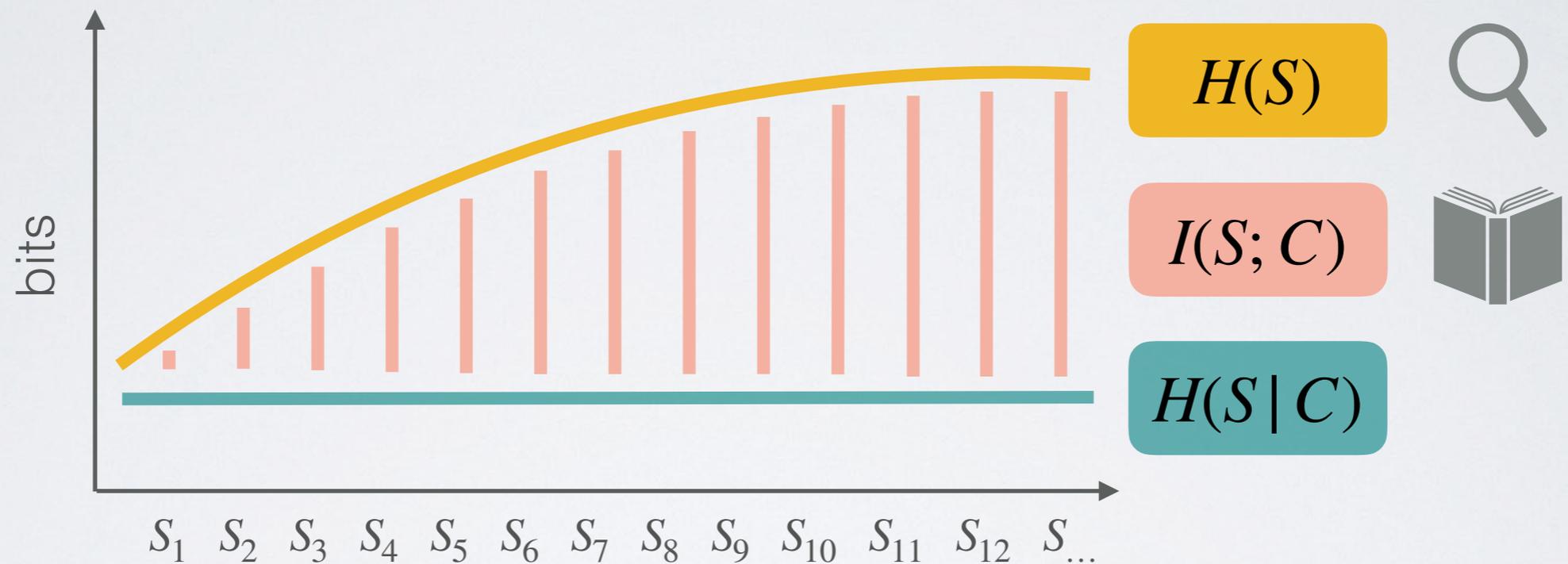
# Constancy Rate Principle

(Genzel & Charniak, 2002)



# Constancy Rate Principle

(Genzel & Charniak, 2002)

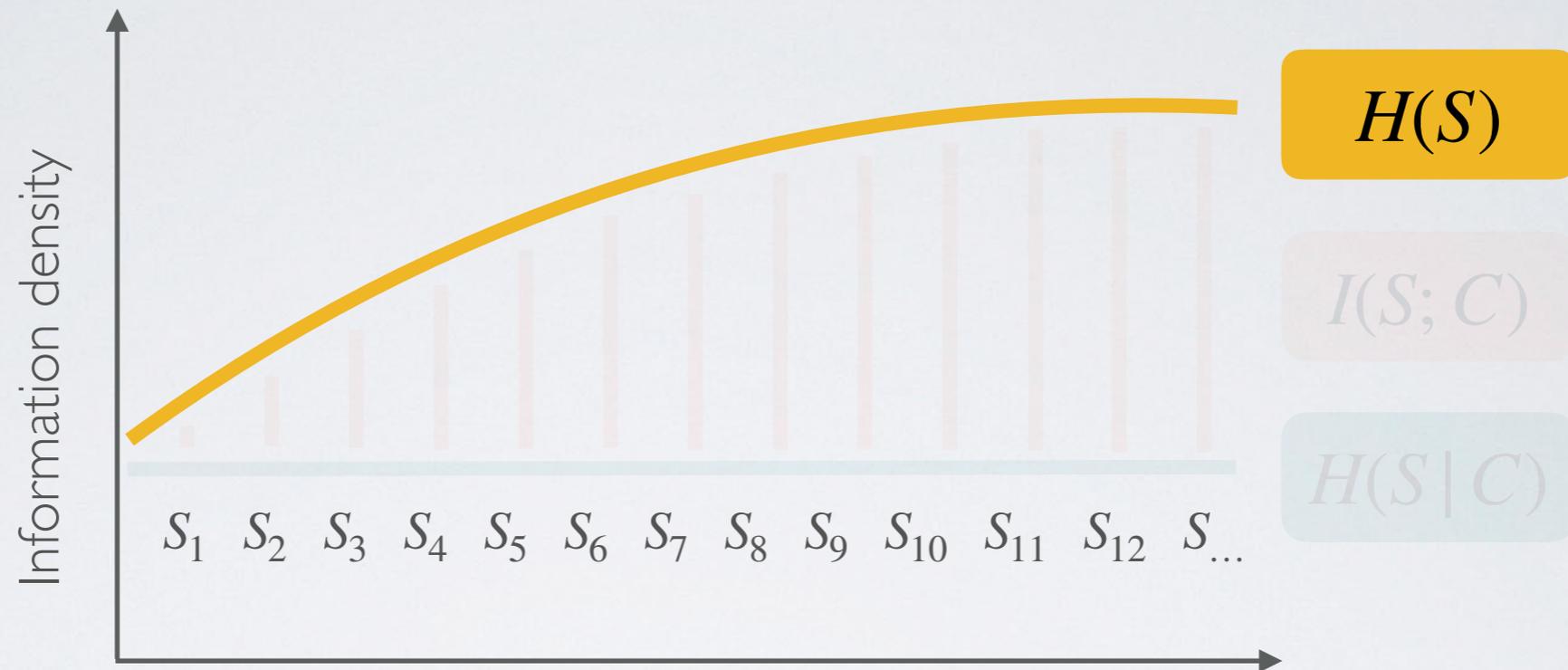


 Assume: as discourse develops, context informativeness increases

 Measure: whether out-of-context surprisal increases

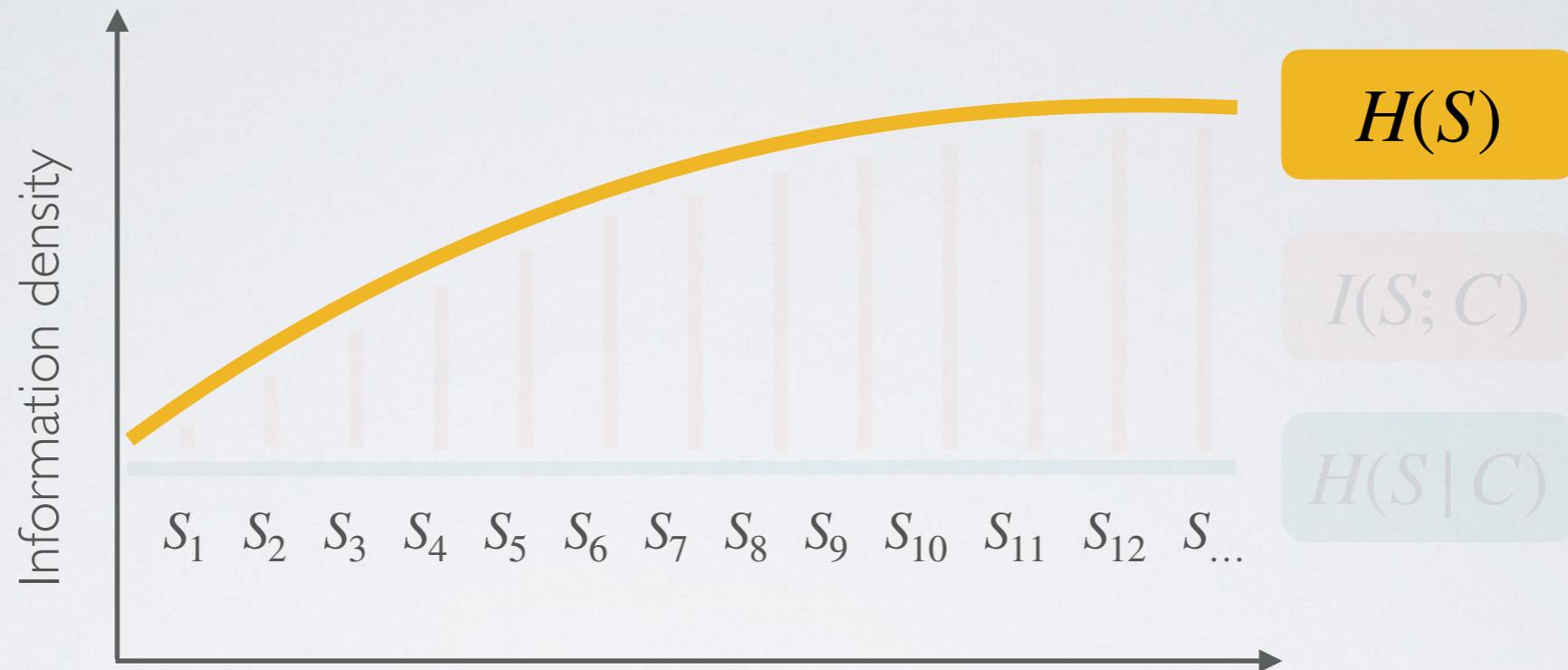
$$H(S | C) \equiv H(S) - I(S; C)$$

# Constancy Rate Principle



We consider the Constancy Rate Principle to hold if  $H(S_i)$  increases with dialogue turn  $i$

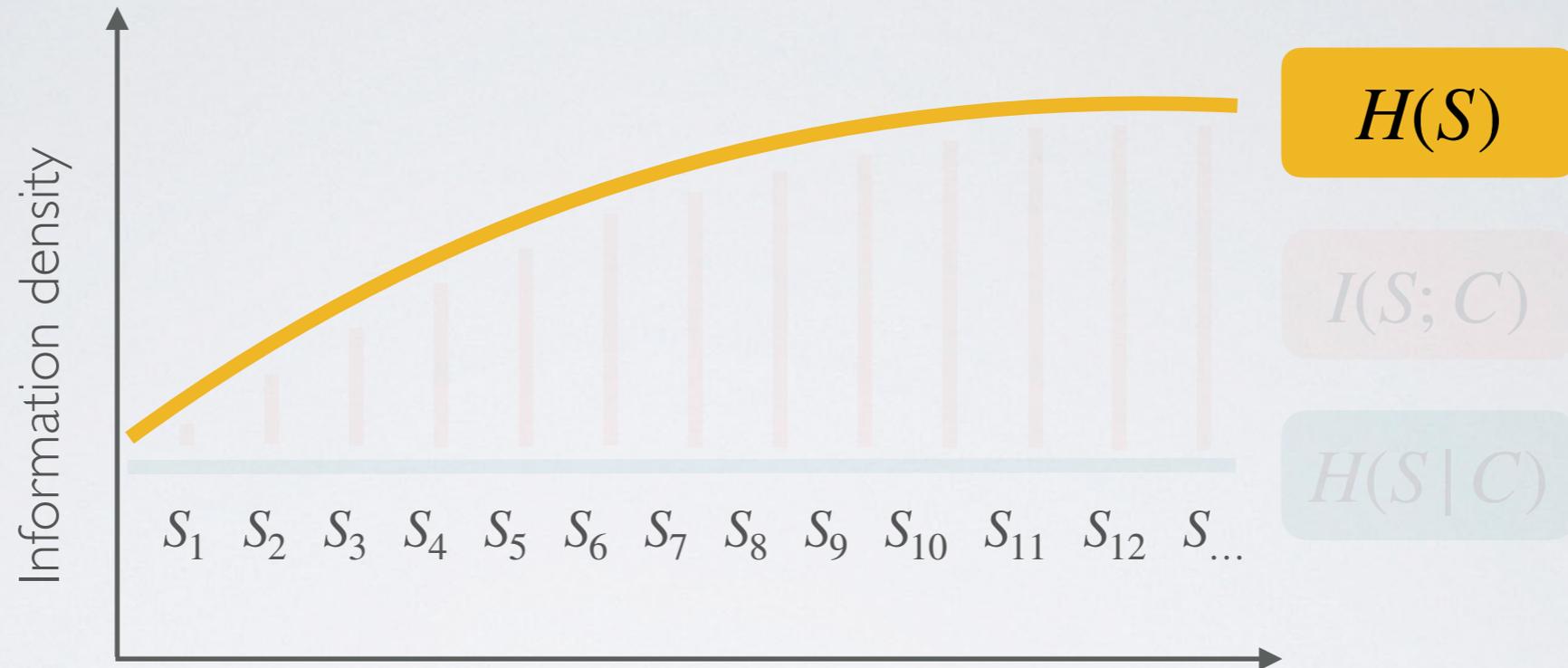
# Estimates of surprisal



$$H(S) = -\log_2 P(S) = -\frac{1}{|S|} \sum_{w_i \in S} \log_2 P(w_i | w_1, \dots, w_{i-1})$$

$P(w_i | \dots)$  estimates obtained with GPT-2  
fine-tuned on 70% of each target corpus  
(30% held-out for analysis)

# Linear Mixed Effect Model



length-normalised  
(Xu & Reitter, 2018)

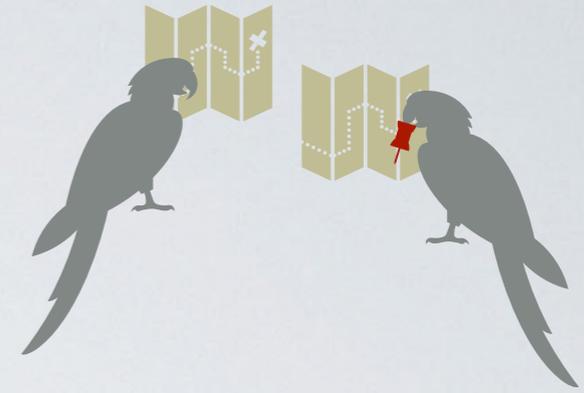
$H(S)$

$$\sim 1 + \log \text{dialogue turn} + (1 + \log \text{dialogue turn} | \text{dialogue\_id})$$

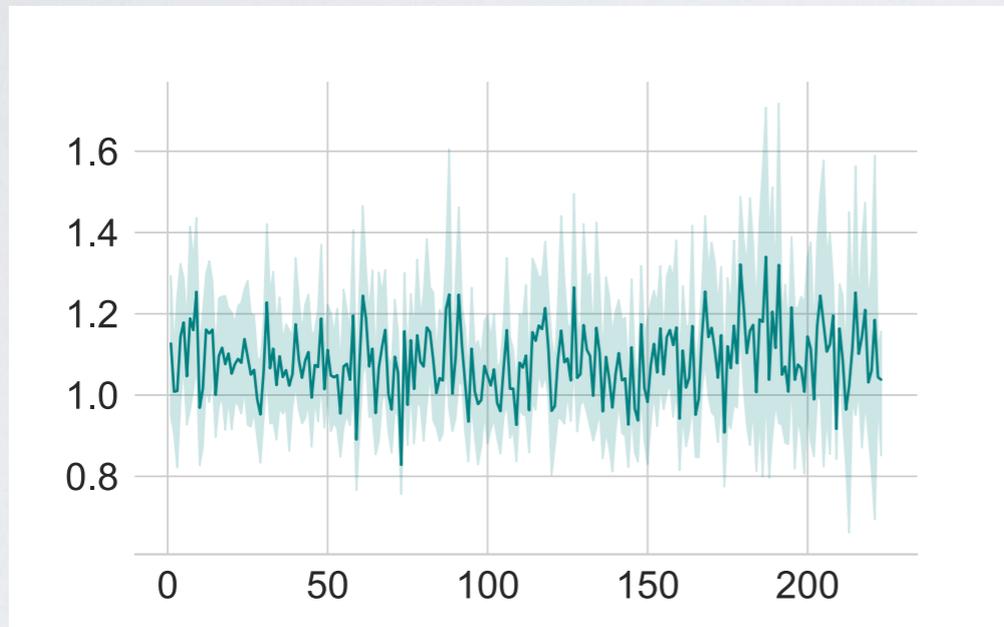
fixed effect

random intercept and slope  
grouped by dyad

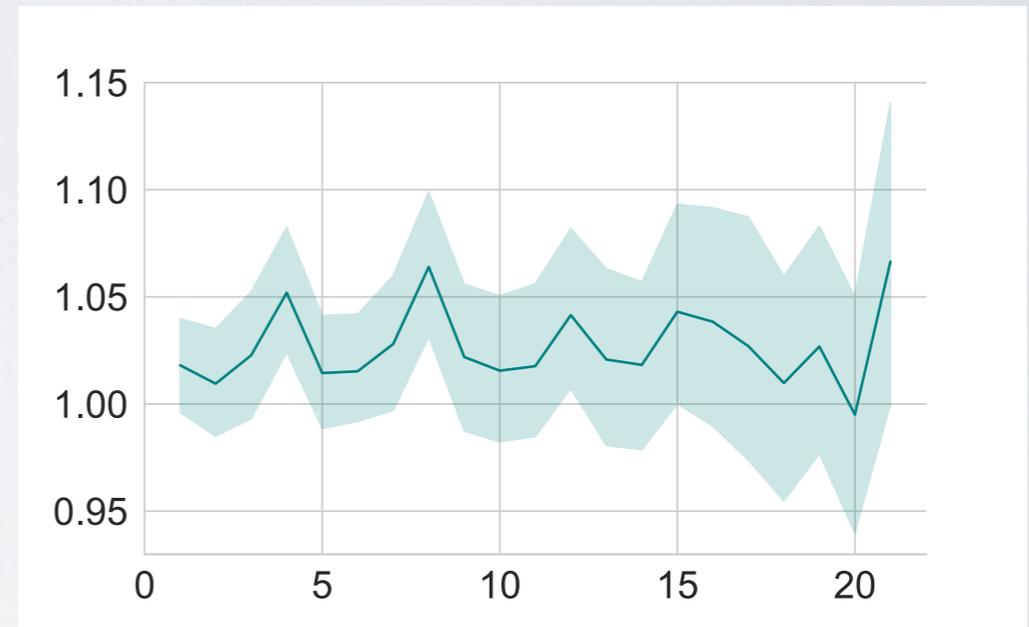
# Results: MapTask



$H(S_i)$



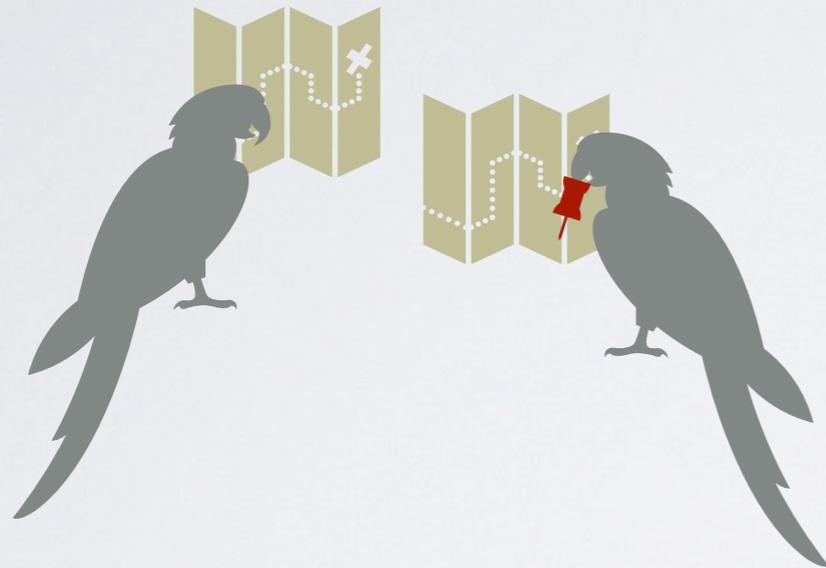
→ dialogue turn  $i$



→ transaction turn  $i$

$H(S)$  does not increase  
over an entire dialogue and over a transaction

# Information transmission acts



**MapTask**  
 instruction giving  
 and following  
 (Anderson et al., 1991)

G	okay
G	go down
F	down towards missionary camp
G	aye
G	just right down uh-huh
F	right
F	i'm i've stopped right above it
G	then
G	you've what
F	i'm right above the missionary camp

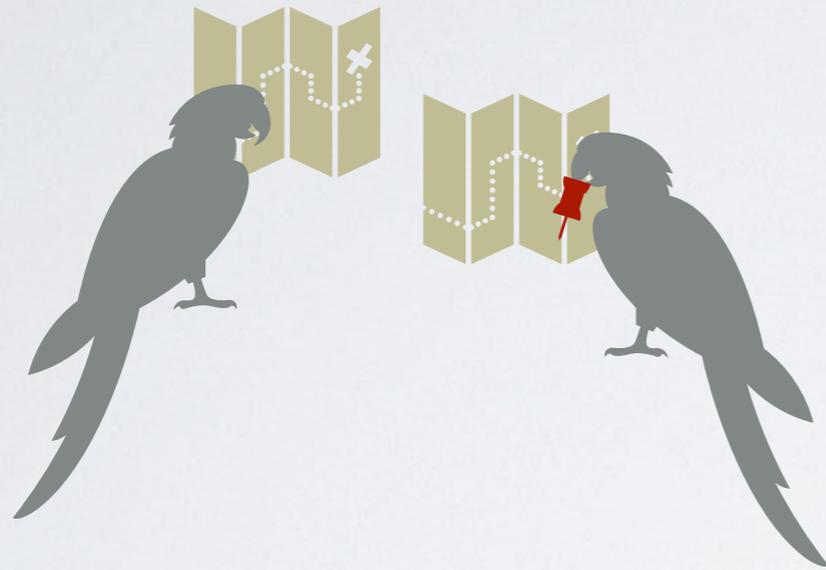
**Transaction 2**

G	just go right down to round down to [...]
F	right
G	okay
F	so go round it
G	uh-huh
G	round to the left uh the left yeah
F	keep going down
G	uh-huh
G	just stop where you when you reach [...]

**Transaction 3**

# Information transmission acts

(no backchannels and grounding acts: *okay, mmhmm*)



**MapTask**  
instruction giving  
and following  
(Anderson et al., 1991)

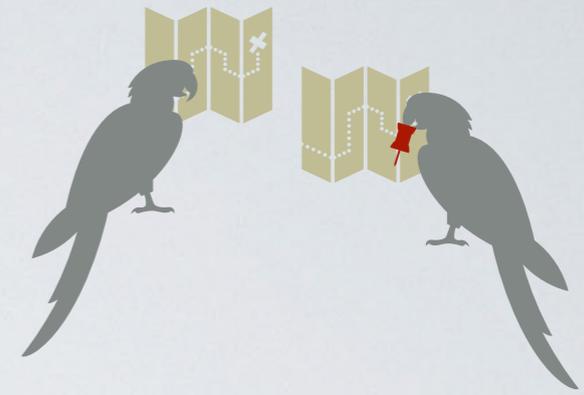
G	okay
G	go down
F	down towards missionary camp
G	aye
G	just right down uh-huh
F	right
F	i'm i've stopped right above it
G	then
G	you've what
F	i'm right above the missionary camp

Transaction 2

G	just go right down to round down to [...]
F	right
G	okay
F	so go round it
G	uh-huh
G	round to the left uh the left yeah
F	keep going down
G	uh-huh
G	just stop where you when you reach [...]

Transaction 3

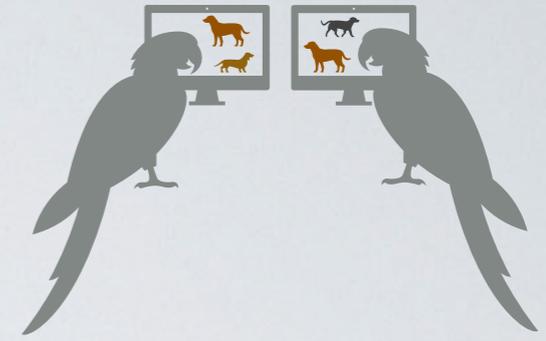
# Results: MapTask



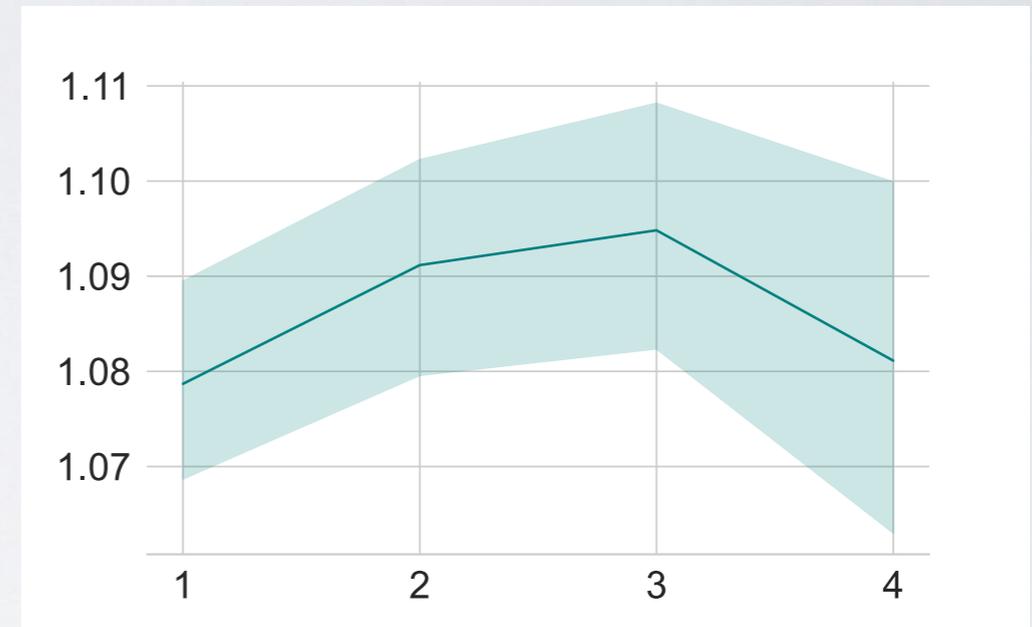
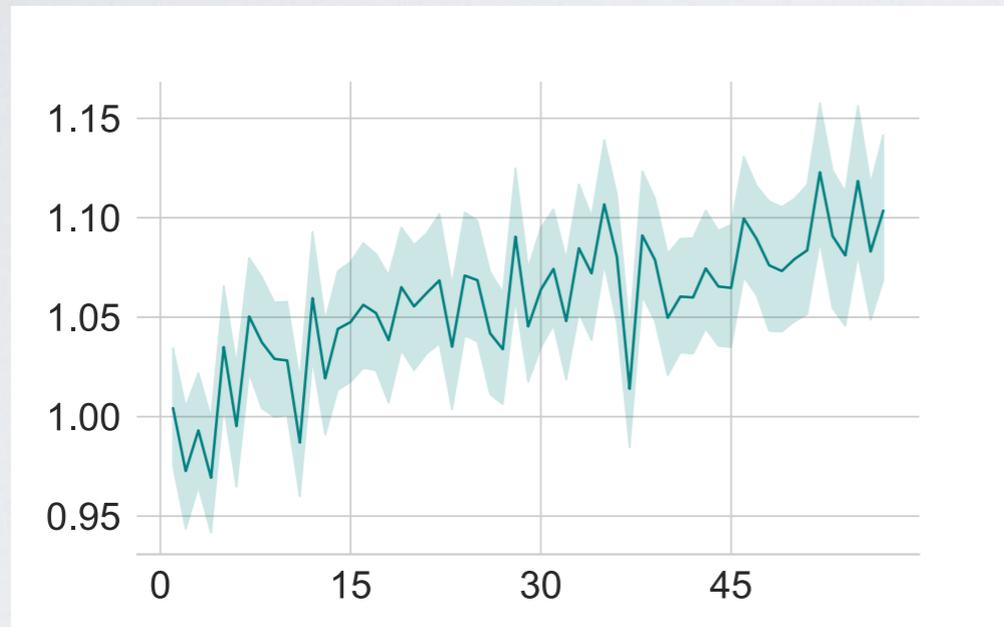
$H(S)$  increases

in the information-transmission dialogue acts of transactions  
(no backchannels and grounding acts: *okay, mmhmm*)

# Results: PhotoBook



$H(S_i)$



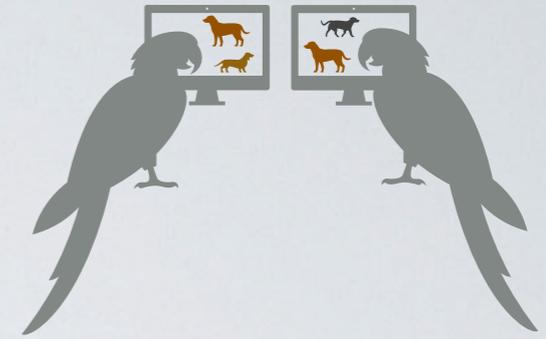
dialogue turn  $i$



reference chain index  $i$

$H(S)$  increases  
over an entire dialogue and over a reference chain

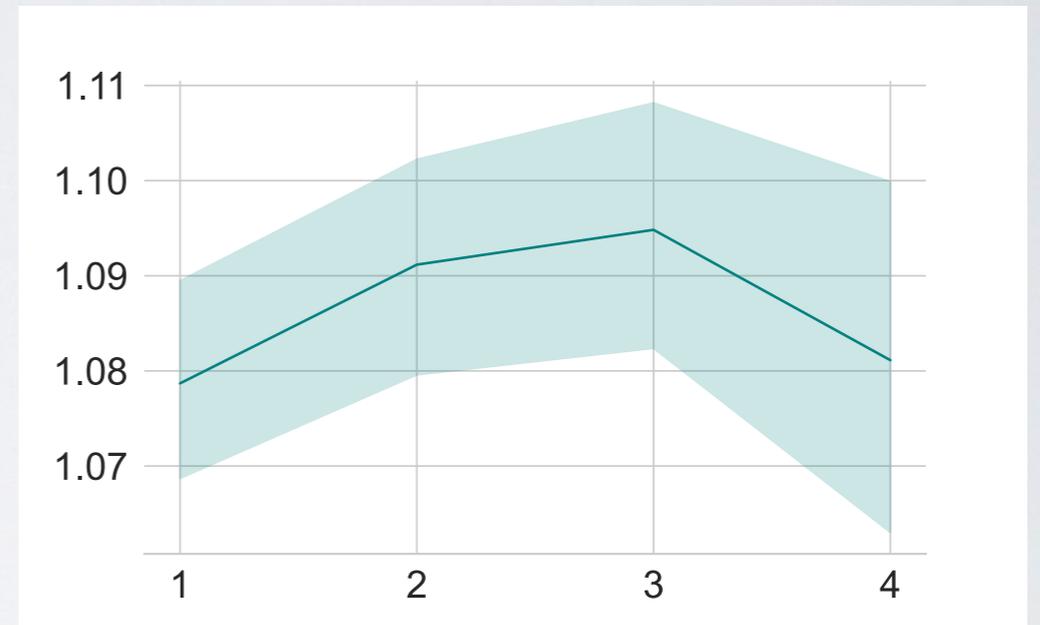
# Results: PhotoBook



$$H(S_i)$$

man eating slice of pizza	0.69
last one for me is guy with pizza	0.78
pizza eater	0.91
pizza	0.67

$$H(S_i)$$



Reduction + information compression

# Study I - Summary

The Constancy Rate Principle holds within more **topically and referentially coherent contextual units** of task-oriented dialogues:

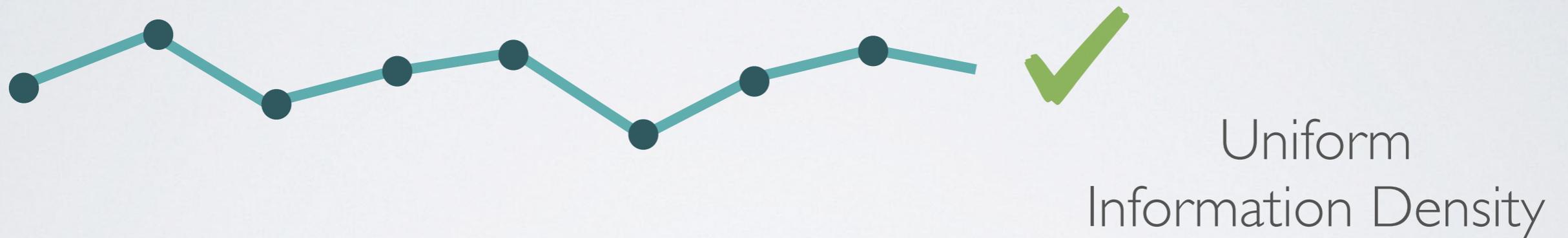
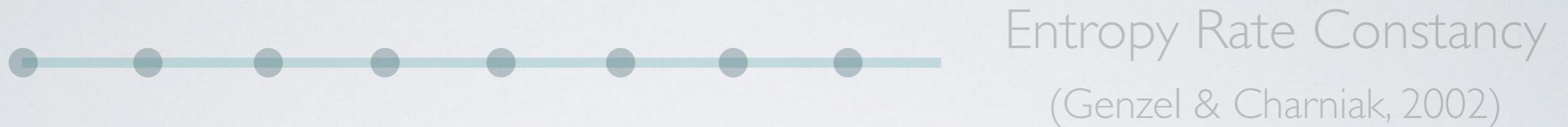
MapTask information transmission acts in transactions  
PhotoBook dialogues and reference chains

# Strategies of information transmission



Constancy Rate Principle  
(Genzel & Charniak, 2002)

# Strategies of information transmission



$S_1$   $S_2$   $S_3$   $S_4$   $S_5$   $S_6$   $S_7$   $S_8$   $S_{\dots}$

## Study 2

**Are surprisal patterns better described as constant (CRP) or as uniform (UID)?**

## Study 2

**Are surprisal patterns better described as constant (CRP) or as uniform (UID)?**

Here, we measure surprisal **as a function of discourse context.**

# Estimates of surprisal (revised)

$$H(S) = -\log_2 P(S) = -\frac{1}{|S|} \sum_{w_i \in S} \log_2 P(w_i | w_1, \dots, w_{i-1})$$

# Estimates of surprisal (revised)

$$H(S) = -\log_2 P(S) = -\frac{1}{|S|} \sum_{w_i \in S} \log_2 P(w_i | w_1, \dots, w_{i-1})$$

$$H(S | C) = -\log_2 P(S | C) = -\frac{1}{|S|} \sum_{w_i \in S} \log_2 P(w_i | w_1, \dots, w_{i-1}, \underline{C})$$

# Estimates of surprisal (revised)

$$H(S) = -\log_2 P(S) = -\frac{1}{|S|} \sum_{w_i \in S} \log_2 P(w_i | w_1, \dots, w_{i-1})$$

$$H(S | C) = -\log_2 P(S | C) = -\frac{1}{|S|} \sum_{w_i \in S} \log_2 P(w_i | w_1, \dots, w_{i-1}, \underline{C})$$

$$I(S; C) \equiv H(S) - H(S | C) = -\log_2 P(S) + \log_2 P(S | C)$$

# Estimates of surprisal (revised)

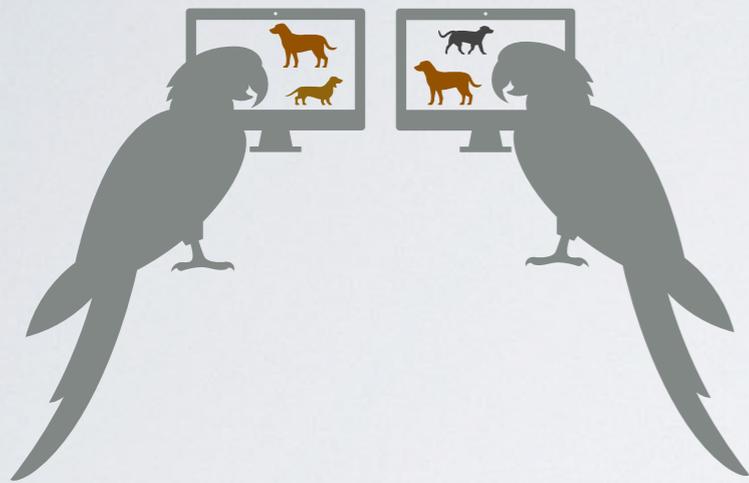
$$H(S) = -\log_2 P(S) = -\frac{1}{|S|} \sum_{w_i \in S} \log_2 P(w_i | w_1, \dots, w_{i-1})$$

$$H(S | C) = -\log_2 P(S | C) = -\frac{1}{|S|} \sum_{w_i \in S} \log_2 P(w_i | w_1, \dots, w_{i-1}, \underline{C})$$

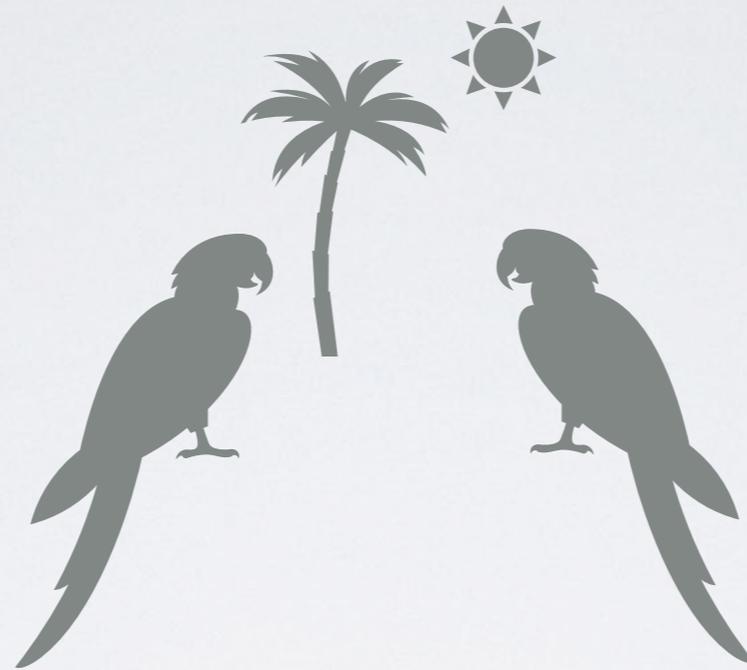
$$I(S; C) \equiv H(S) - H(S | C) = -\log_2 P(S) + \log_2 P(S | C)$$

$P(w_i | \dots)$  estimates obtained with GPT-2  
fine-tuned on 70% of each target corpus  
(30% held-out for analysis)

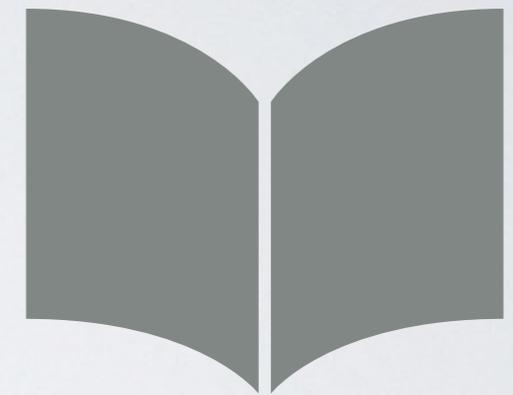
# Experimental data



**PhotoBook**  
task-oriented  
written dialogues  
(Haber et al., 2019)

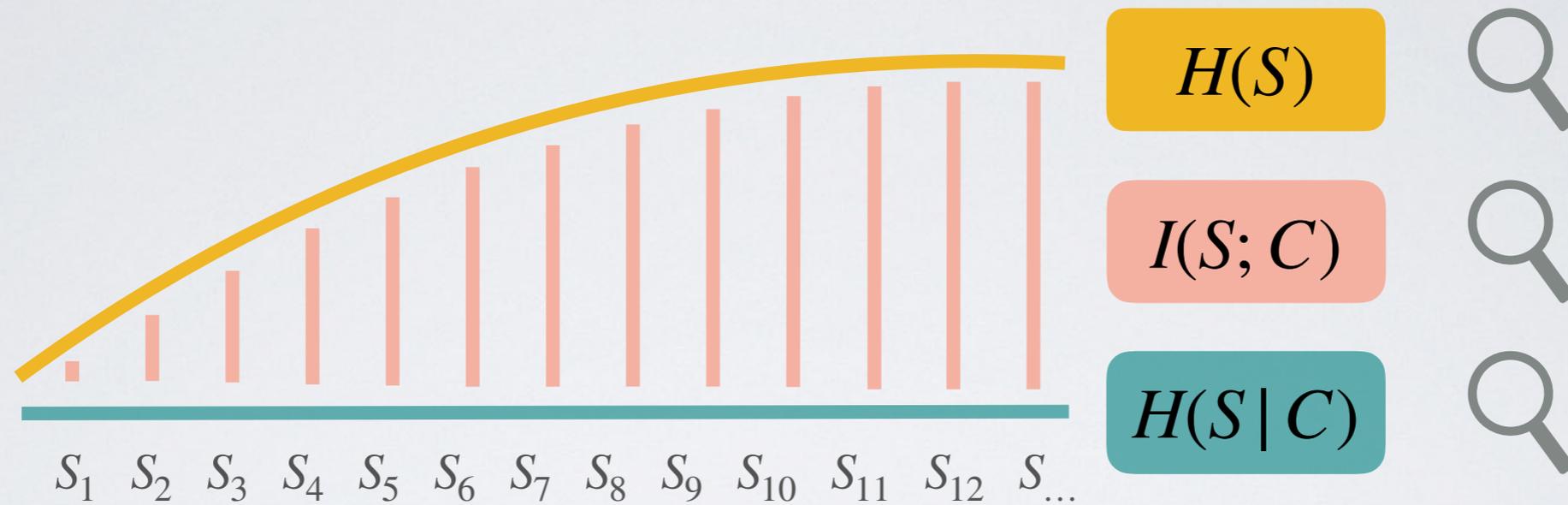


**Spoken BNC**  
open domain  
dialogues  
(Love et al., 2017)

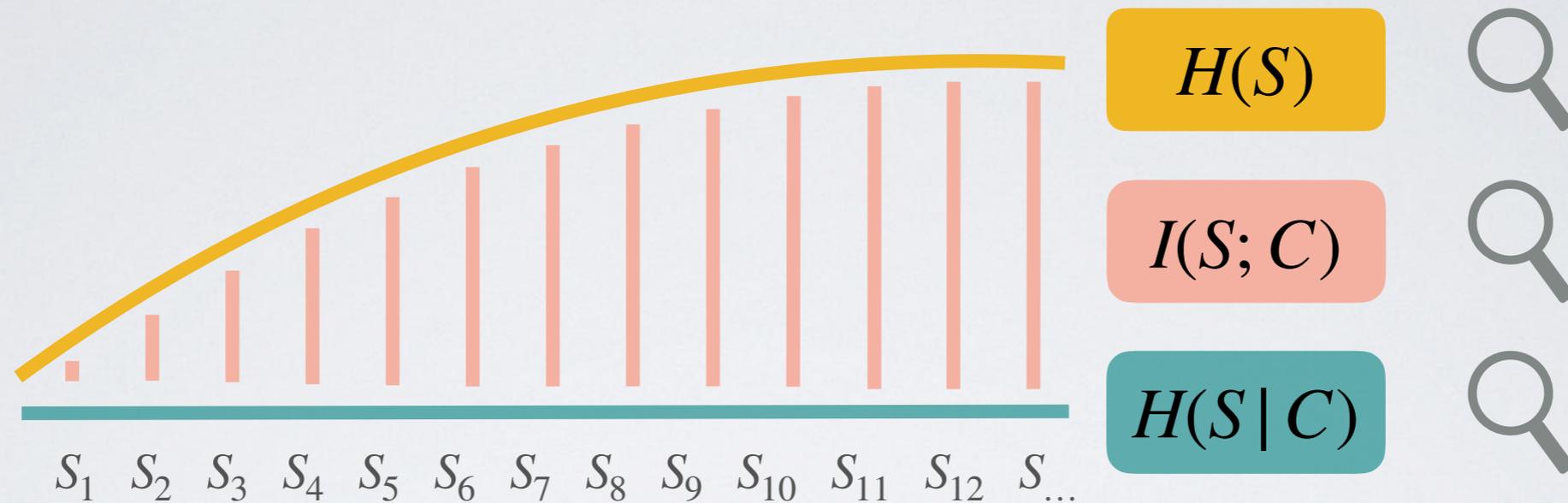


**Penn Treebank**  
newspaper  
articles  
(Mitchell et al. 1999)

# Can surprisal be described as constant?

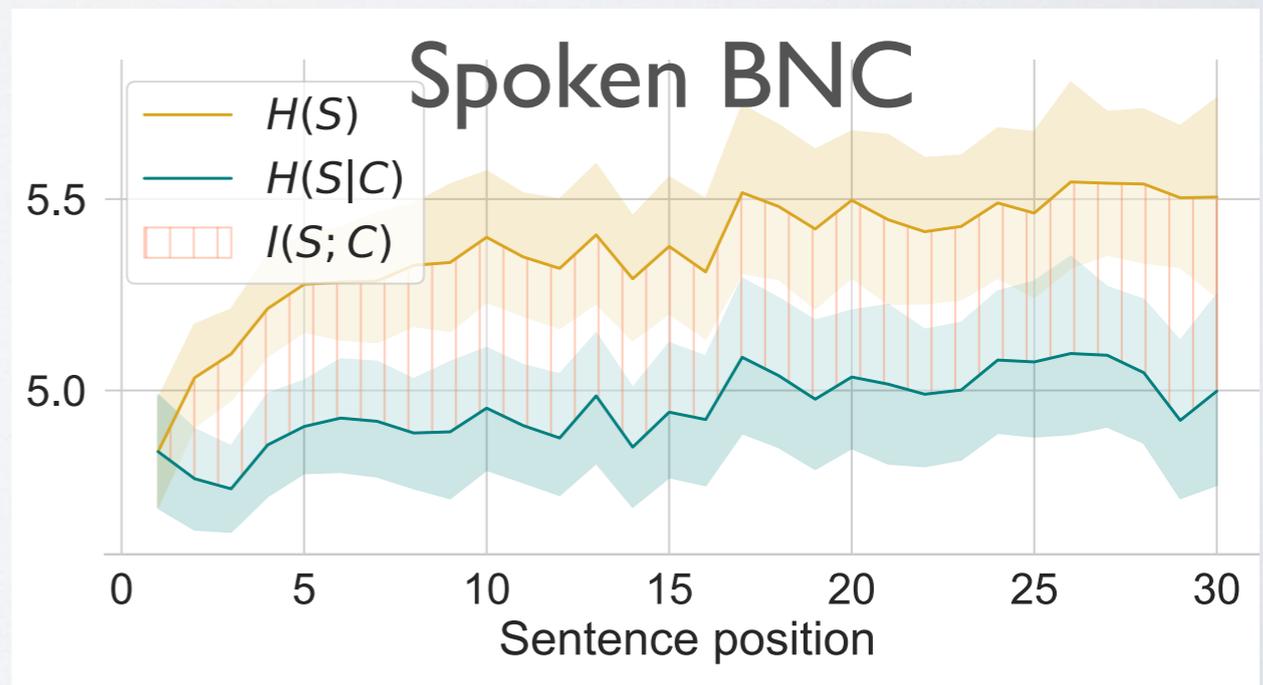
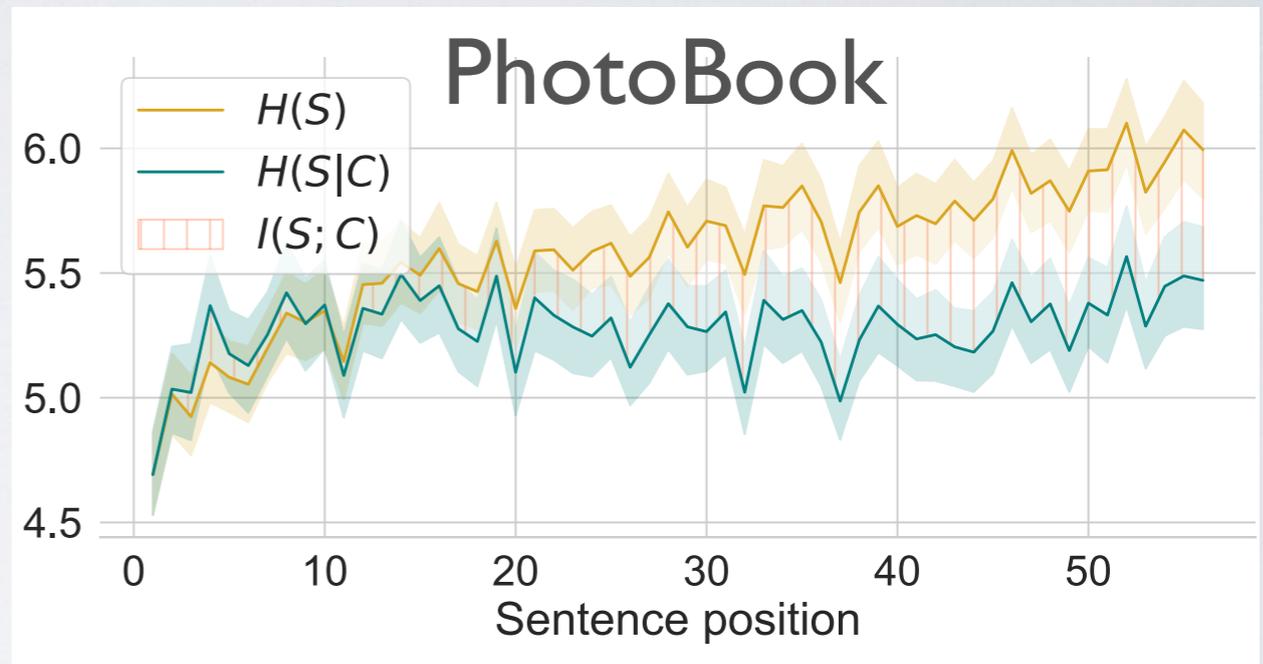
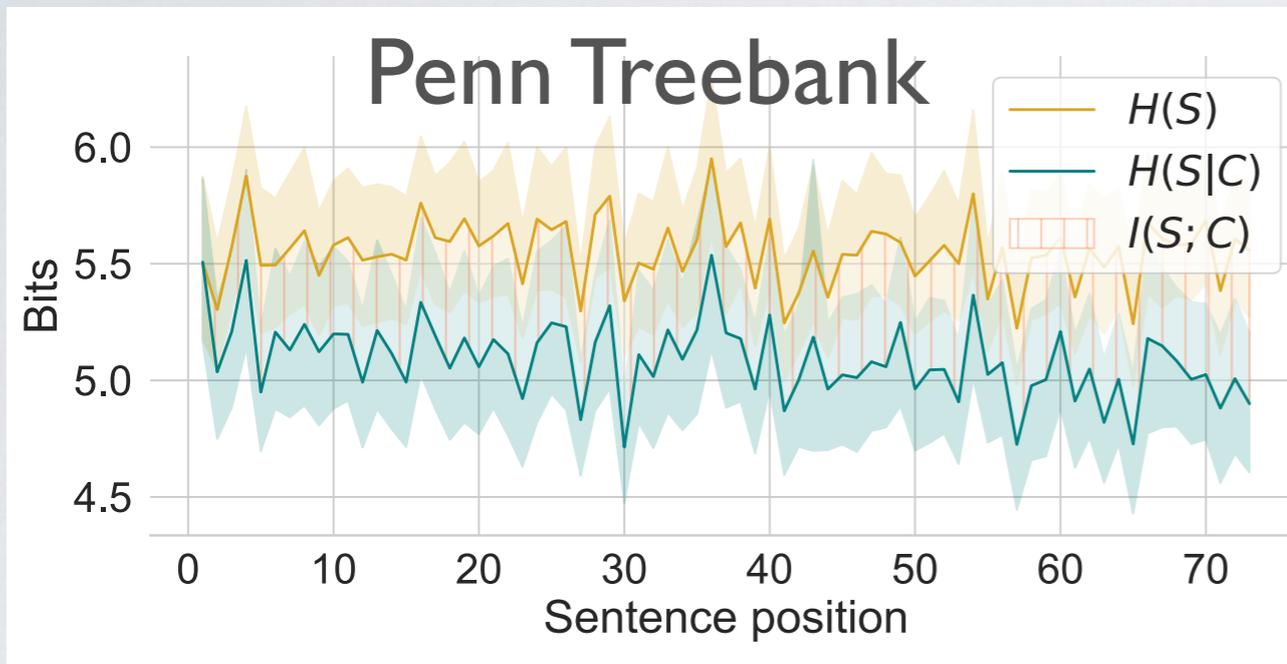
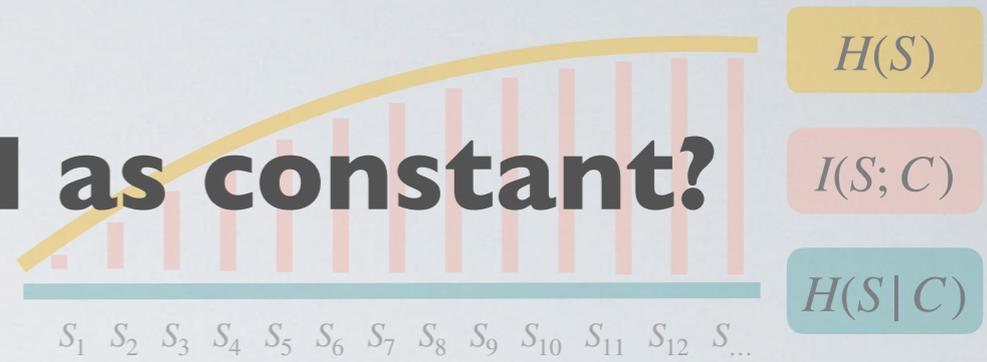


# Can surprisal be described as constant?

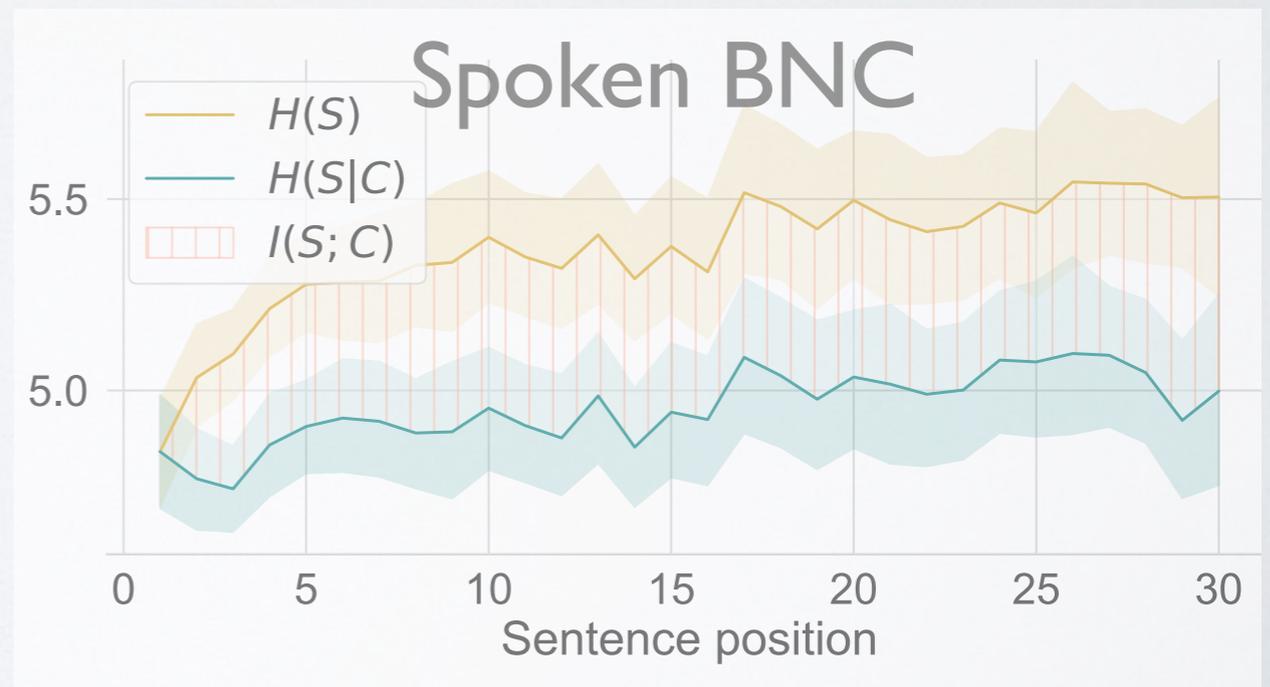
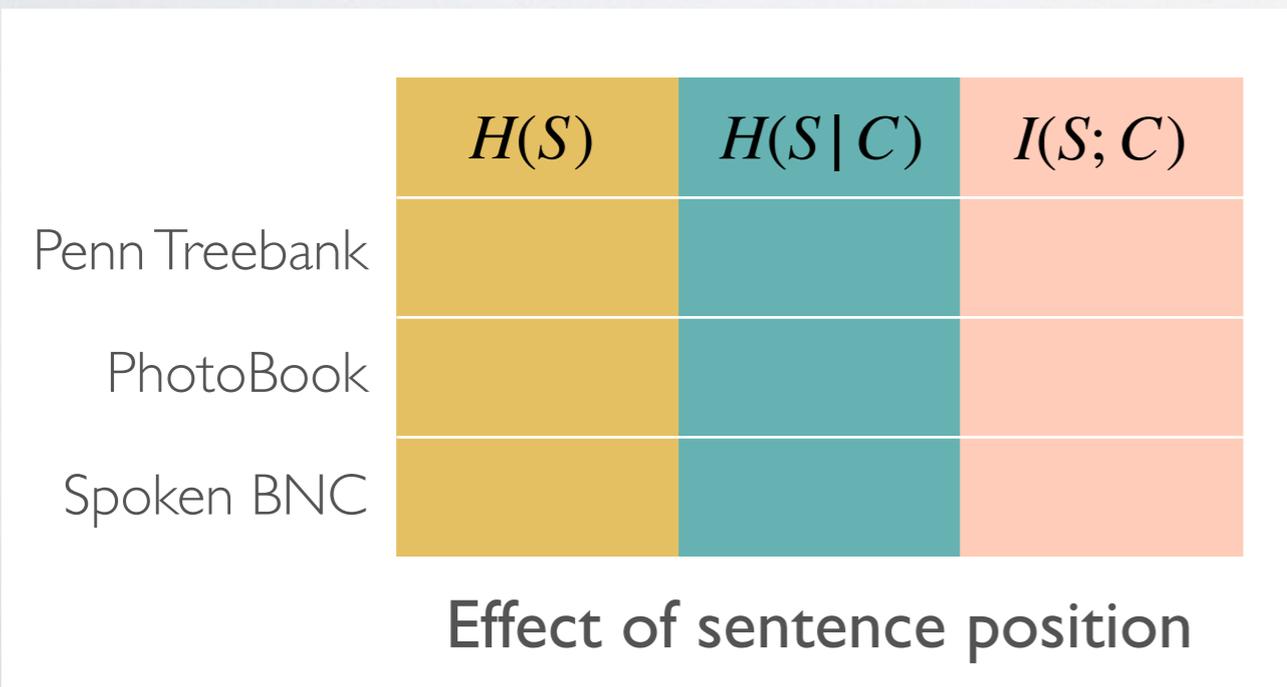
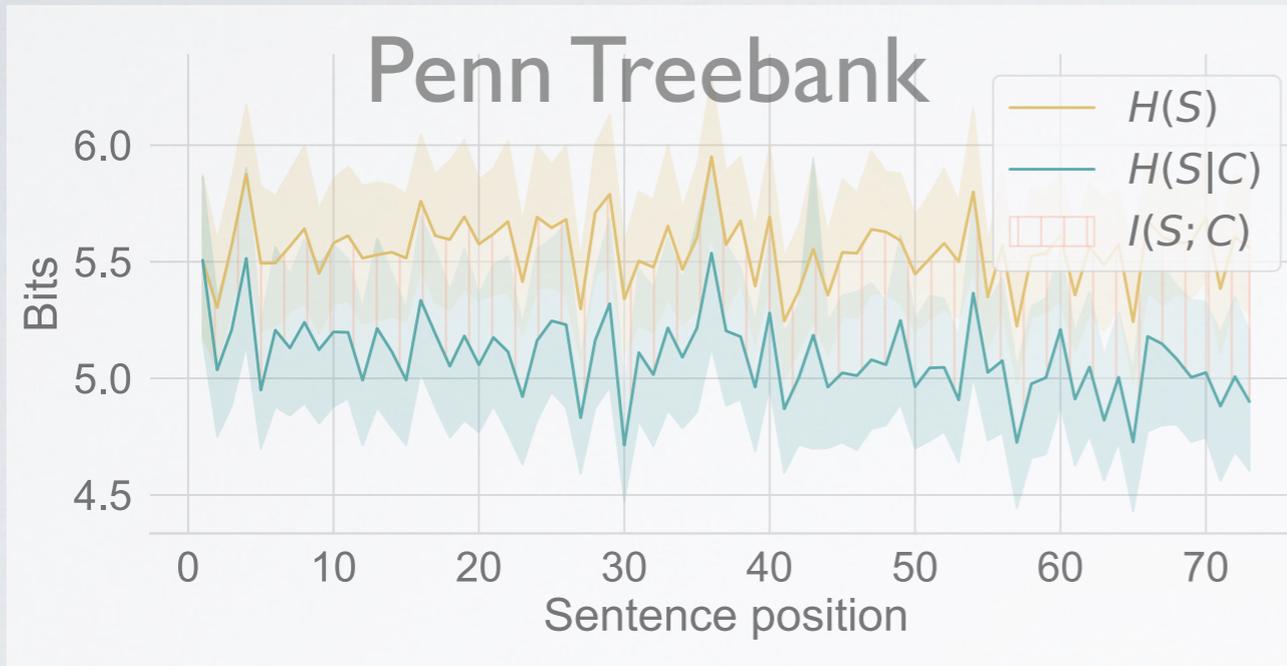


$H(S)$   $\sim$  fixed effects  
 $I(S; C)$   $\sim$   $1 + \log \text{ sentence position} + \log \text{ sentence length}$   
 $H(S|C)$   $\sim$   $+ (1 + \log \text{ sentence position} + \log \text{ sentence length} \mid \text{doc\_id})$   
random intercept and slope  
 grouped by text / dialogue

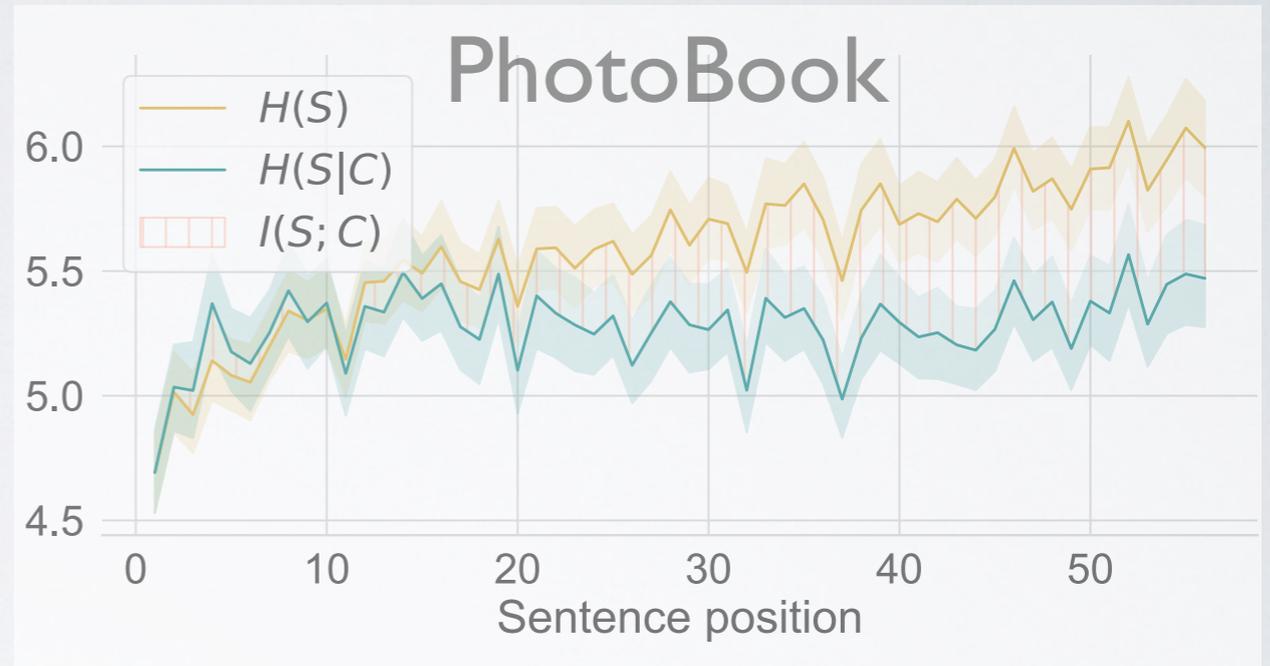
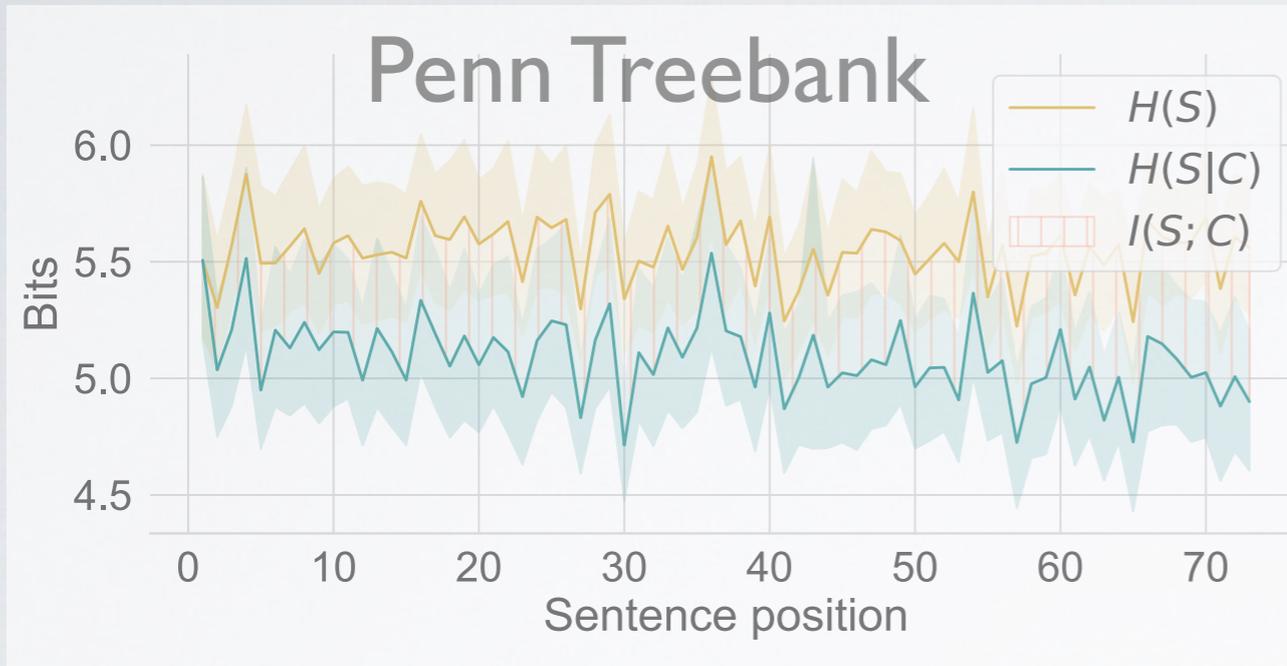
# Can surprisal be described as constant?



# Can surprisal be described as constant?

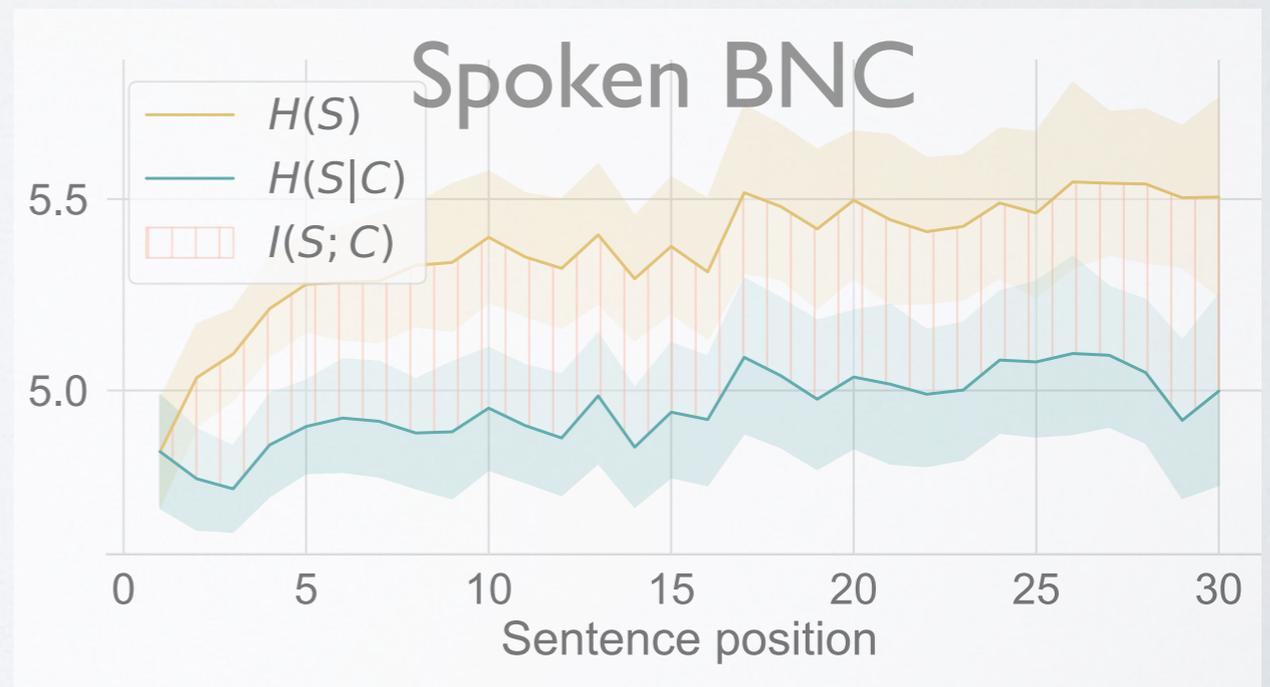


# Can surprisal be described as constant?

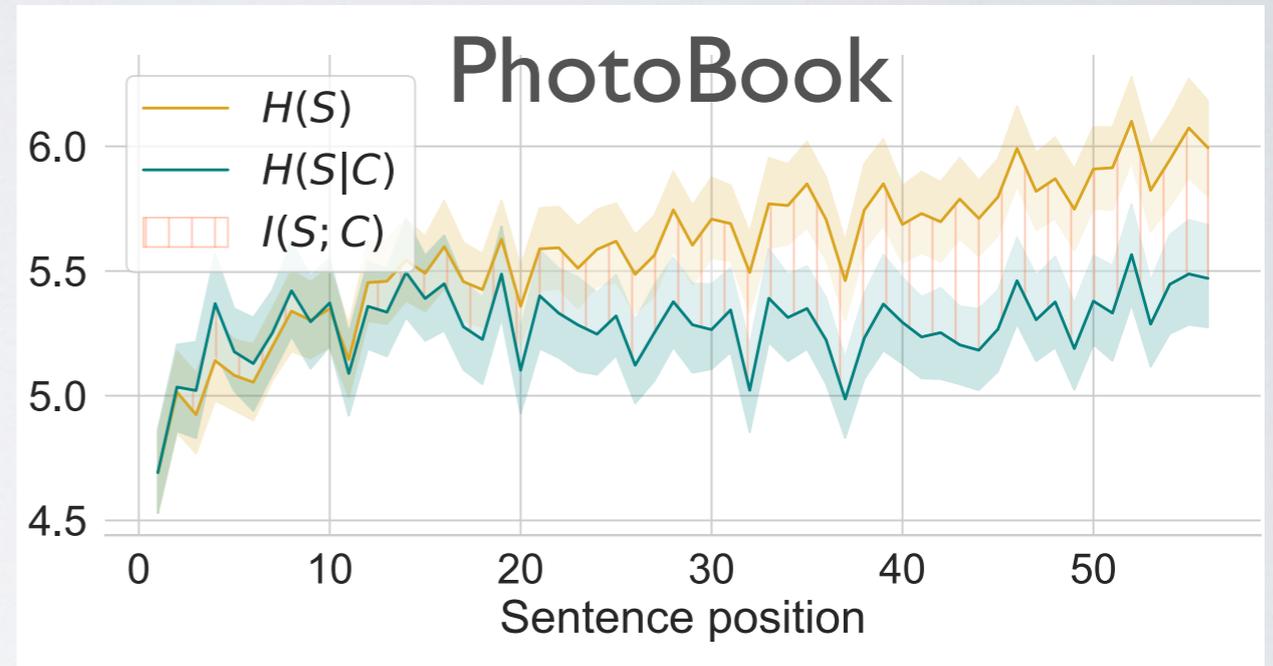
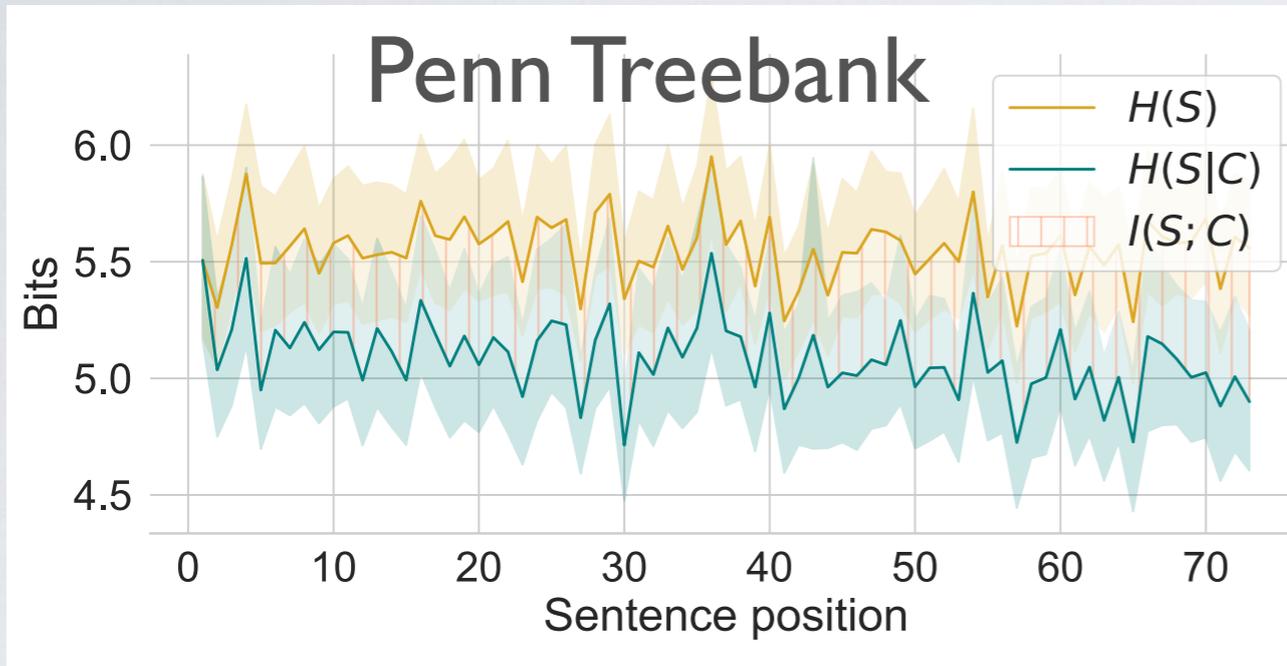


	$H(S)$	$H(S C)$	$I(S; C)$
Penn Treebank			↗
PhotoBook			↗
Spoken BNC			↗

Effect of sentence position

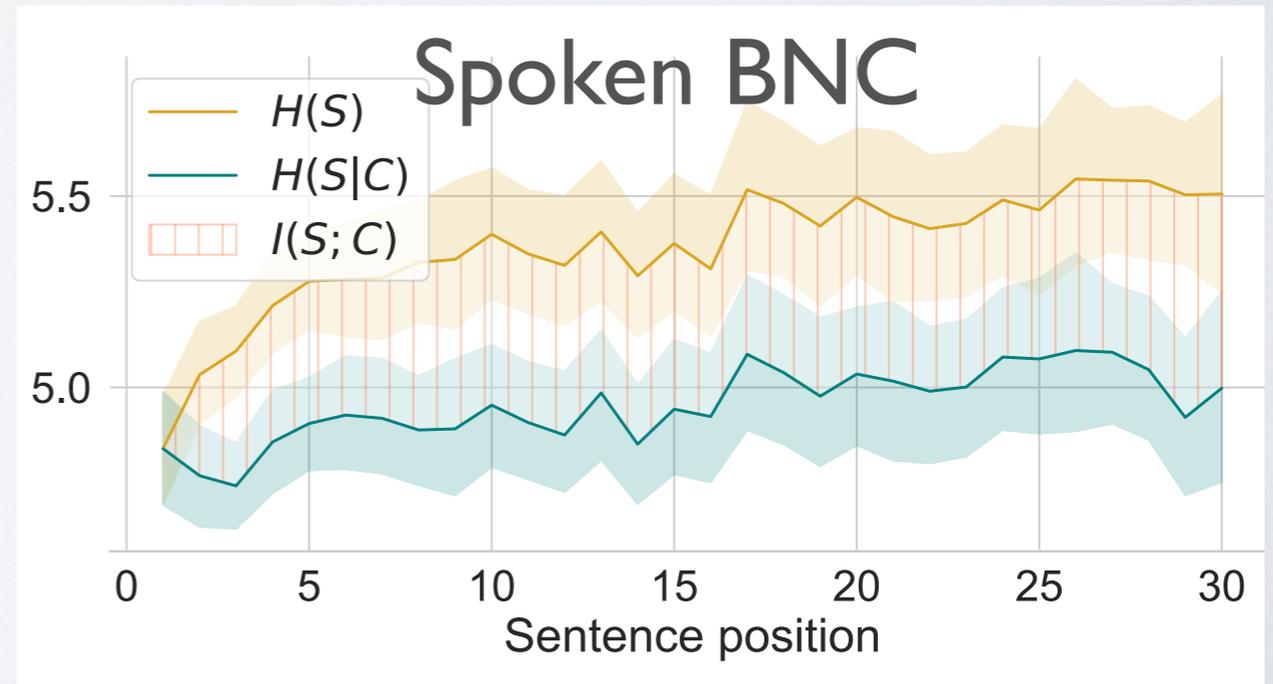


# Can surprisal be described as constant?

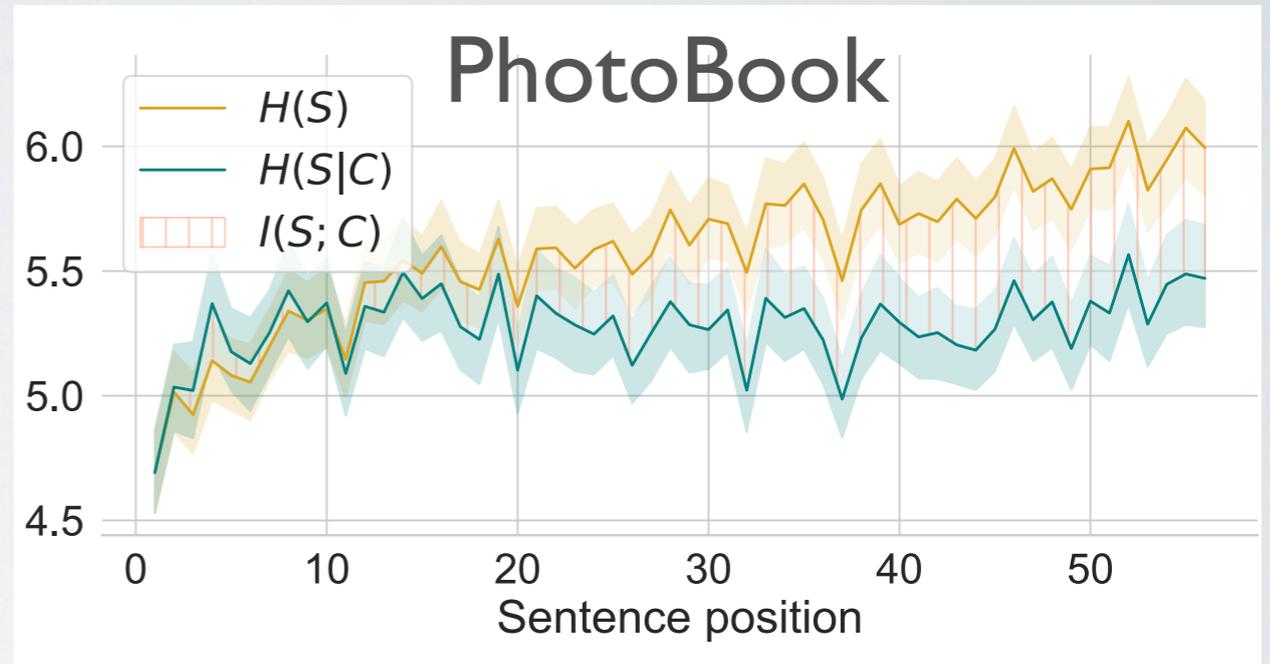
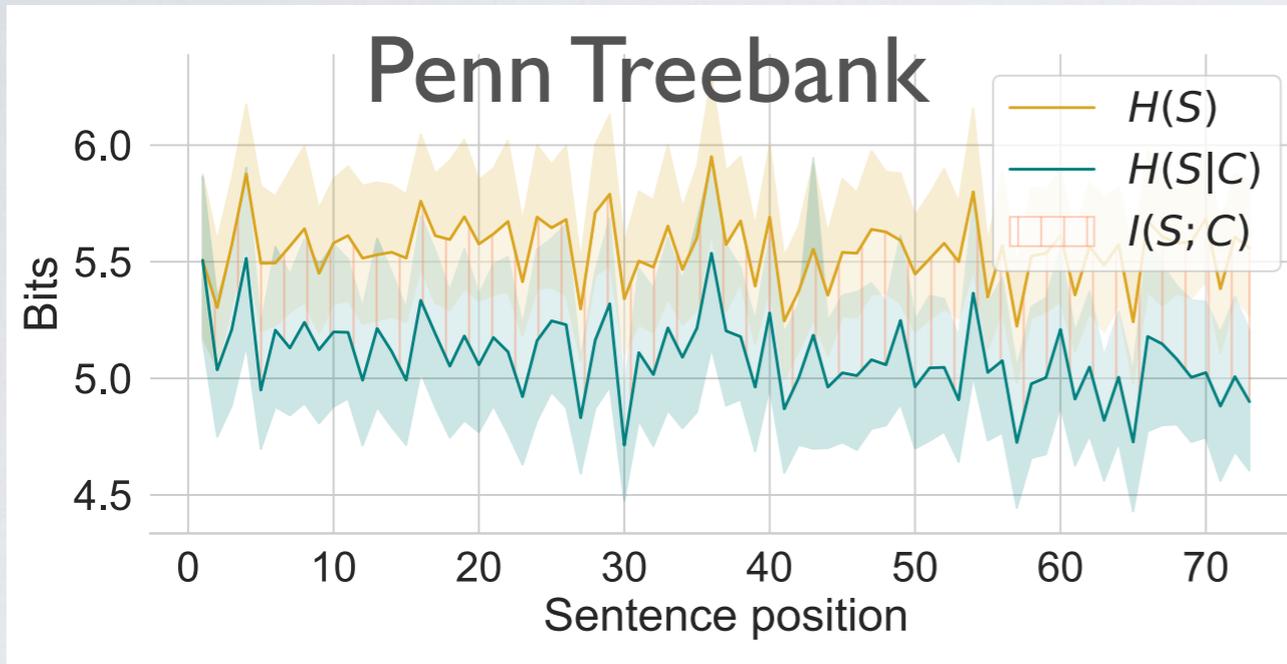
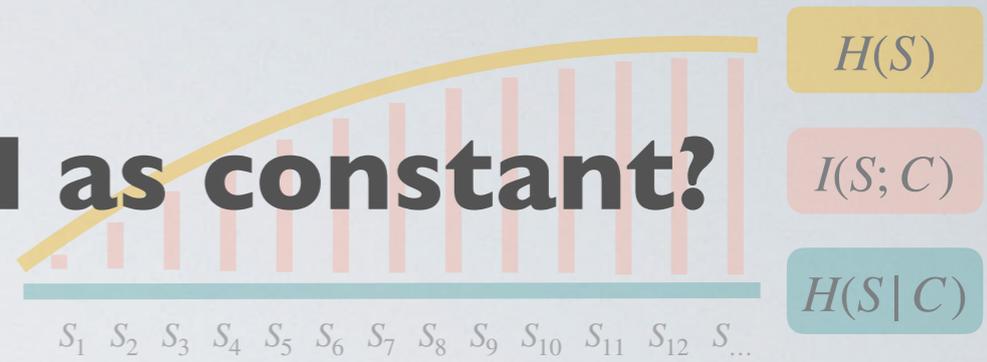


	$H(S)$	$H(S C)$	$I(S; C)$
Penn Treebank	↗		↗
PhotoBook	↗		↗
Spoken BNC	→		↗

Effect of sentence position

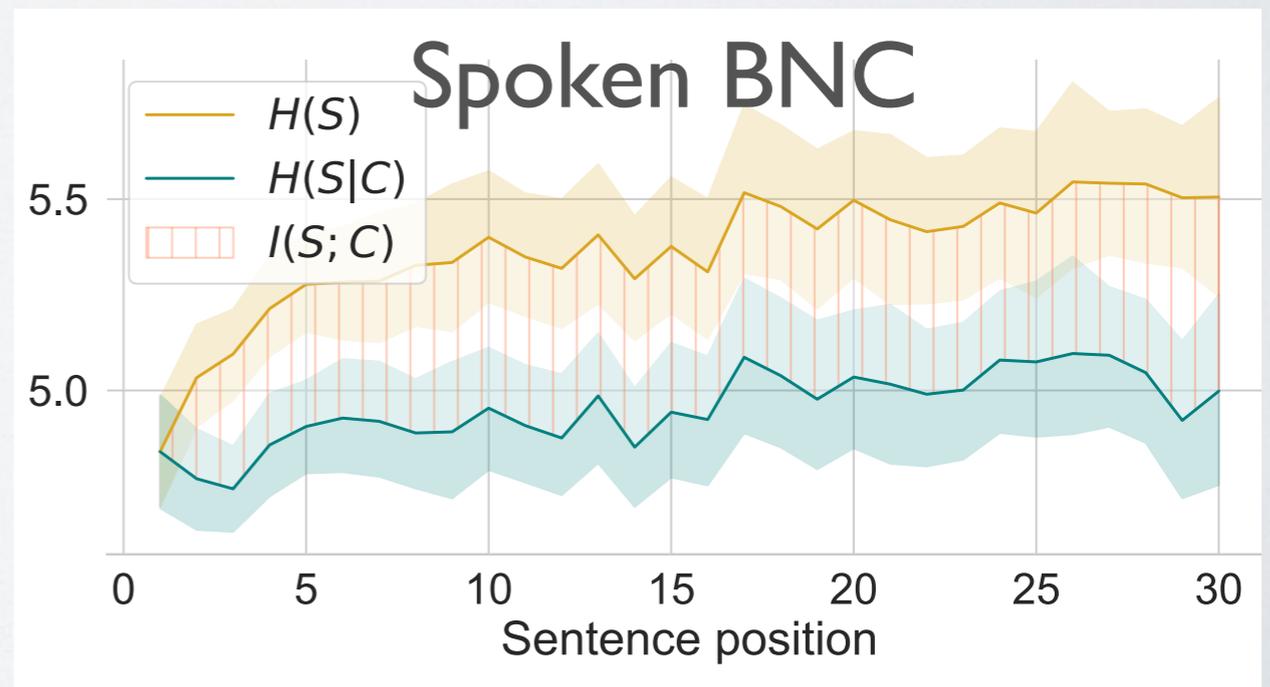


# Can surprisal be described as constant?



	$H(S)$	$H(S C)$	$I(S; C)$
Penn Treebank	↗	↗ (green)	↗
PhotoBook	↗	↘ (red)	↗
Spoken BNC	↗	↘ (red)	↗

Effect of sentence position



**Can surprisal be described as uniform?**

# Can surprisal be described as uniform?

Criteria of uniformity (Collins, 2014)

Global centrality

$$-\frac{1}{N} \sum_{i=1}^N \left( H(S_i | C_i) - \mu \right)^2$$

Local predictability

$$-\frac{1}{N} \sum_{i=2}^N \left( H(S_i | C_i) - H(S_{i-1} | C_{i-1}) \right)^2$$

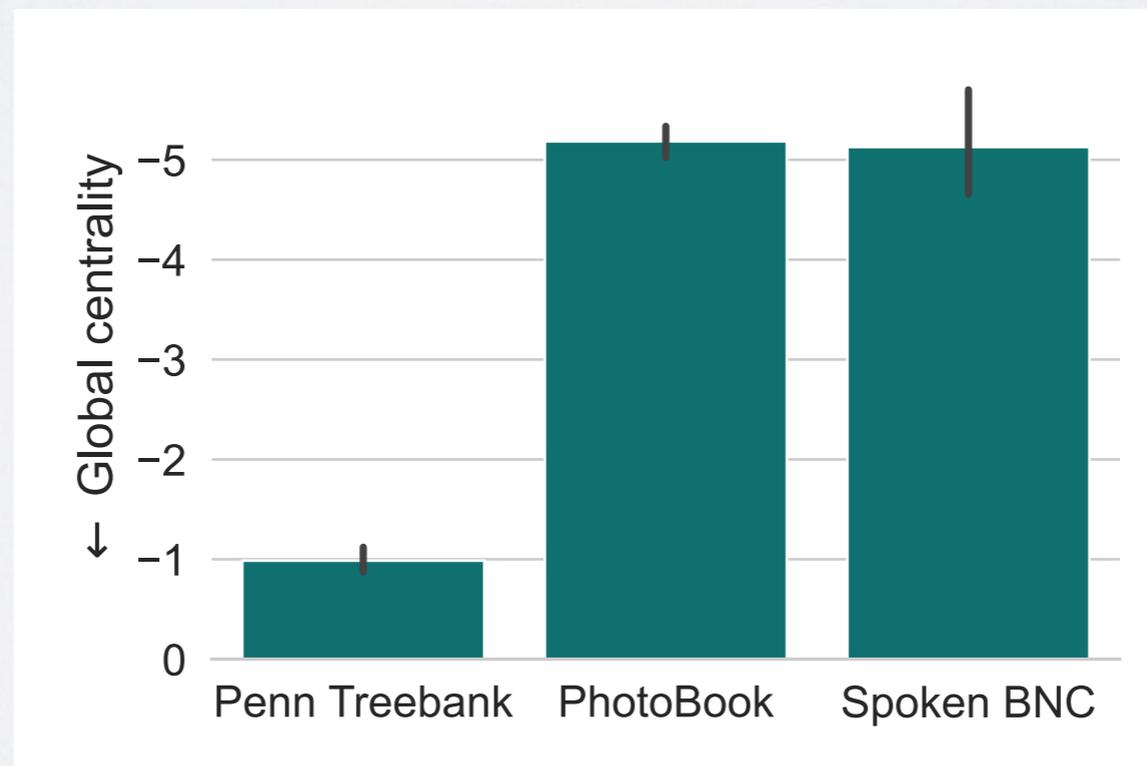
$N$  number of sentences in the text / dialogue

$\mu$  average information content in the text / dialogue

# Can surprisal be described as uniform?

Global centrality

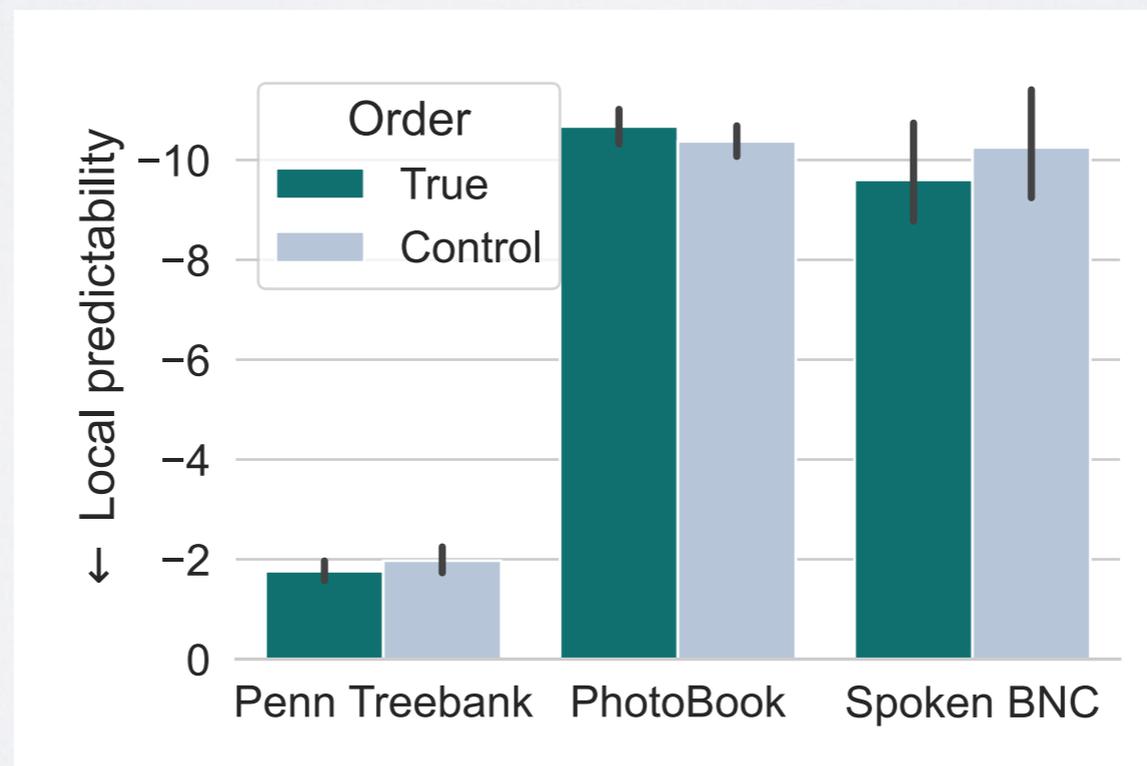
$$-\frac{1}{N} \sum_{i=1}^N \left( H(S_i | C_i) - \mu \right)^2$$



# Can surprisal be described as uniform?

Local predictability

$$-\frac{1}{N} \sum_{i=2}^N \left( H(S_i | C_i) - H(S_{i-1} | C_{i-1}) \right)^2$$



## Study 2 - Summary

We have used Transformers language models to obtain surprisal estimates for sentences **within their discourse context**.

In newspaper articles surprisal remains stable, as predicted by the Constancy Rate Principle.

**Surprisal decreases in dialogues:**  
spoken open domain and written task-oriented dialogues.

**Global uniformity** is a more faithful criterion than local uniformity.

In dialogues, surprisal is stable within  
topically and referentially coherent contextual units...

yet it decreases overall, throughout the entire dialogue

In dialogues, surprisal is stable within  
topically and referentially coherent contextual units...

yet it decreases overall, throughout the entire dialogue

**Any other efficient  
production strategies at play?**

## Study 3

**Does construction repetition reduce surprisal?  
Is it an efficient strategy?**

# Wrapping up

Information-theoretic model of communication to describe decision making in language production.

**Surprisal-based measures** of collaborative effort, estimated via an autoregressive neural language model.

Humans follow **efficient strategies**, leading to near-optimal collaborative effort:

- ★ surprisal decreases throughout dialogues (unlike in written texts)
- ★ topically and referentially coherent contextual units show stable (globally uniform) levels of surprisal
- ★ repetitions of non-topical and non-referential expressions (constructions) reduces surprisal

# What's next?

Apply this analysis to machine-generated utterances, to **evaluate** adherence to human efficiency principles.

Engineer better **training and decoding objectives** for NLG systems, to bring their production strategies closer to those followed by humans.

# Open questions

How to transform information-theoretic measures into evaluation metrics?

How to incorporate them into the training and decoding objectives of NLG systems?

How to improve our model of communication?

- continual adaptation (without catastrophic forgetting)
- partner-specificity
- separate production and processing effort
- include other factors (lexical access, conversational goals, ...)

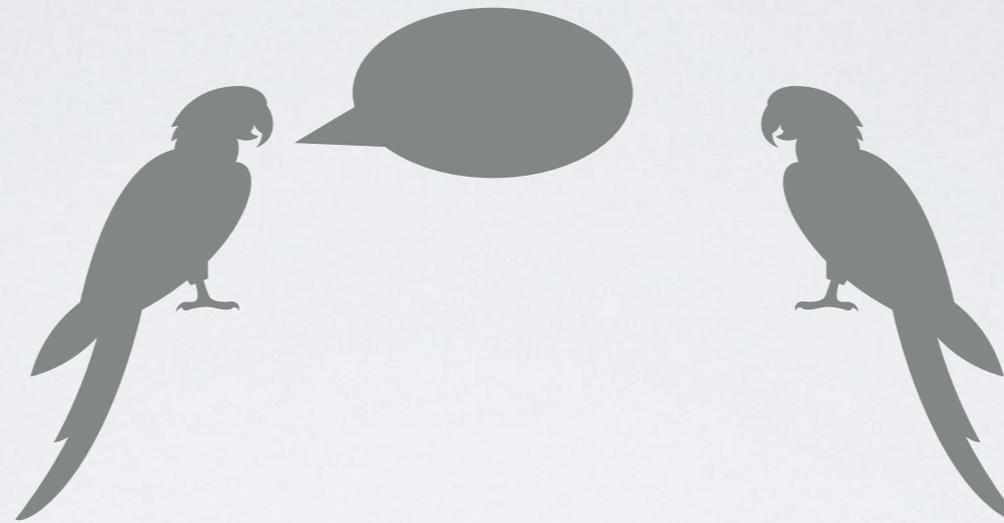
Can we find evidence for efficiency principles in other languages?

# APPENDIX

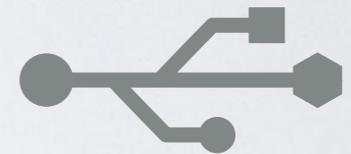
# Strategies of language production

- Grice's maxims and the cooperative principle
- Neo-Gricean principles (Horn and Ward, Levinson)
- Relevance theory (Sperber and Wilson)
- Conversation analysis (Sacks, Schegloff, Jefferson)
- Least collaborative effort (Clark)
- Rational Speech Act framework (Frank and Goodman)

# Collaborative effort in language production



PROCESSING  
EFFORT



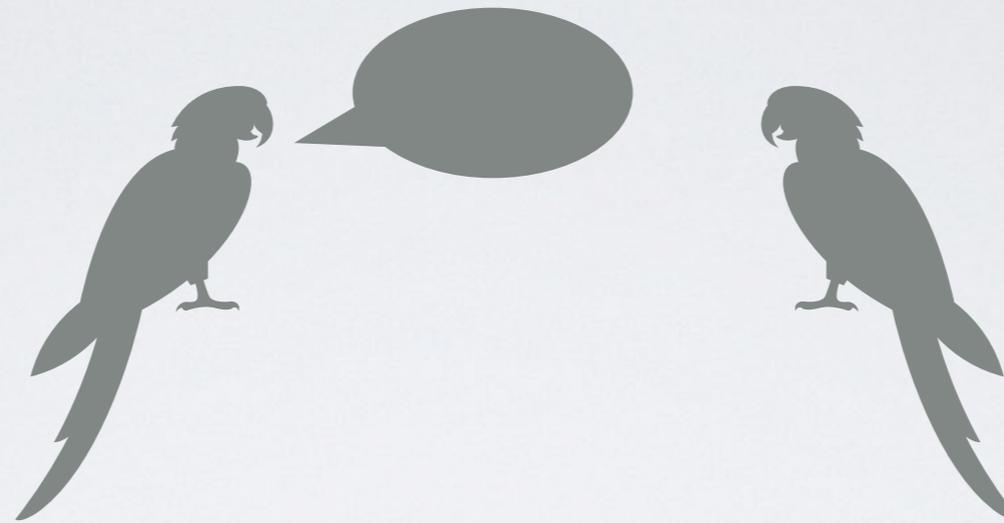
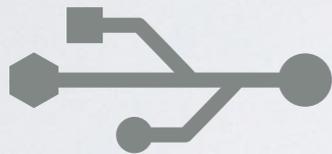
Amcore Financial Inc. said it agreed to acquire Central of Illinois Inc. in a stock swap.

Shareholders of Central, a bank holding company based in Sterling, will receive Amcore stock equal to 10 times Central's 1989 earnings, Amcore said.

**For the first nine months of 1989, Central of Illinois Inc., a bank holding company based in Sterling, earned \$2 million.**

# Collaborative effort in language production

PRODUCTION  
EFFORT



Amcore Financial Inc. said it agreed to acquire Central of Illinois Inc. in a stock swap.

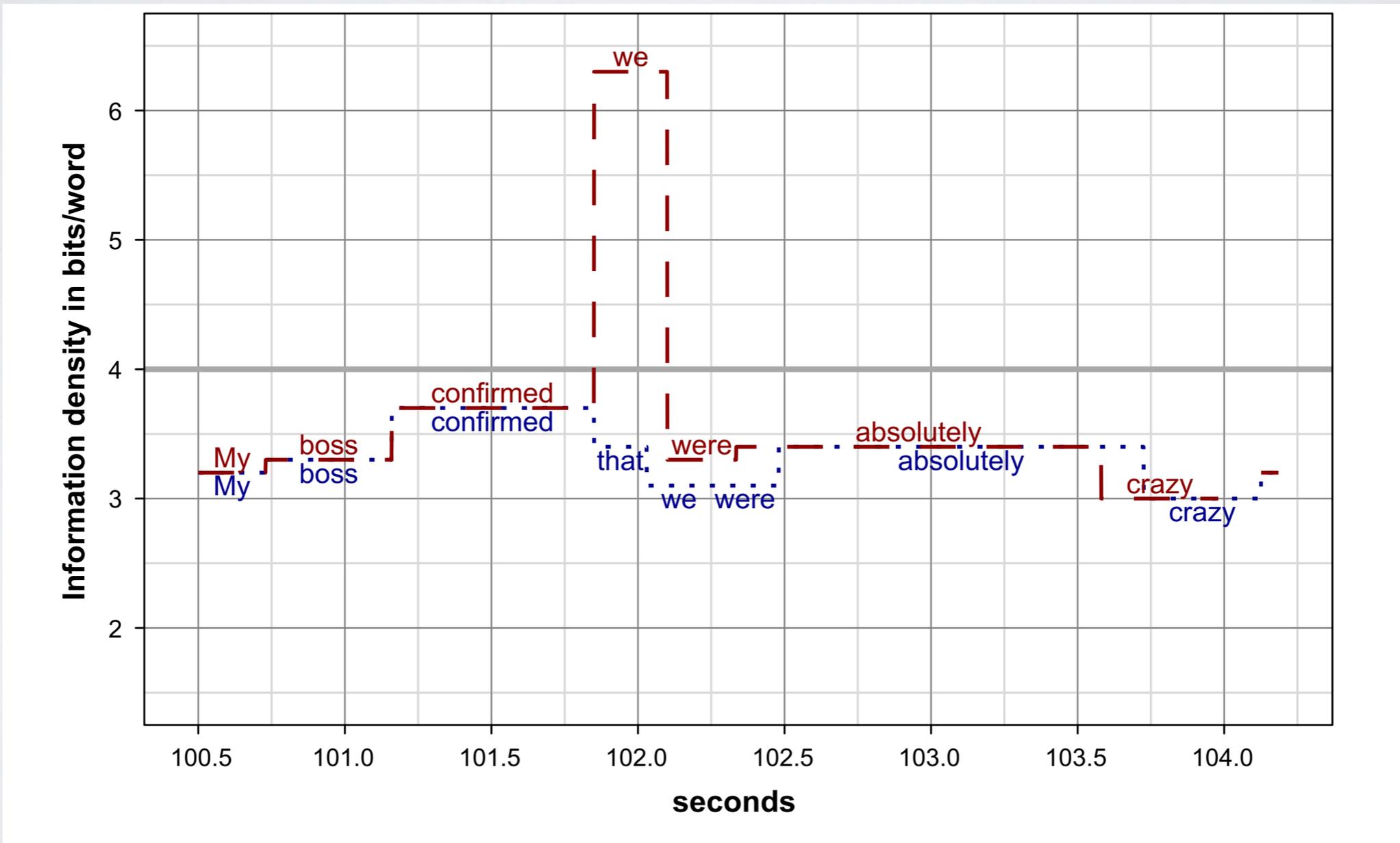
Shareholders of Central, a bank holding company based in Sterling, will receive Amcore stock equal to 10 times Central's 1989 earnings, Amcore said.

**For the first nine months of 1989, Central earned \$2 million.**

# Uniform Information Density

## Example I: Syntax

Complementiser *that*-mentioning (e.g., Jaeger, 2010)

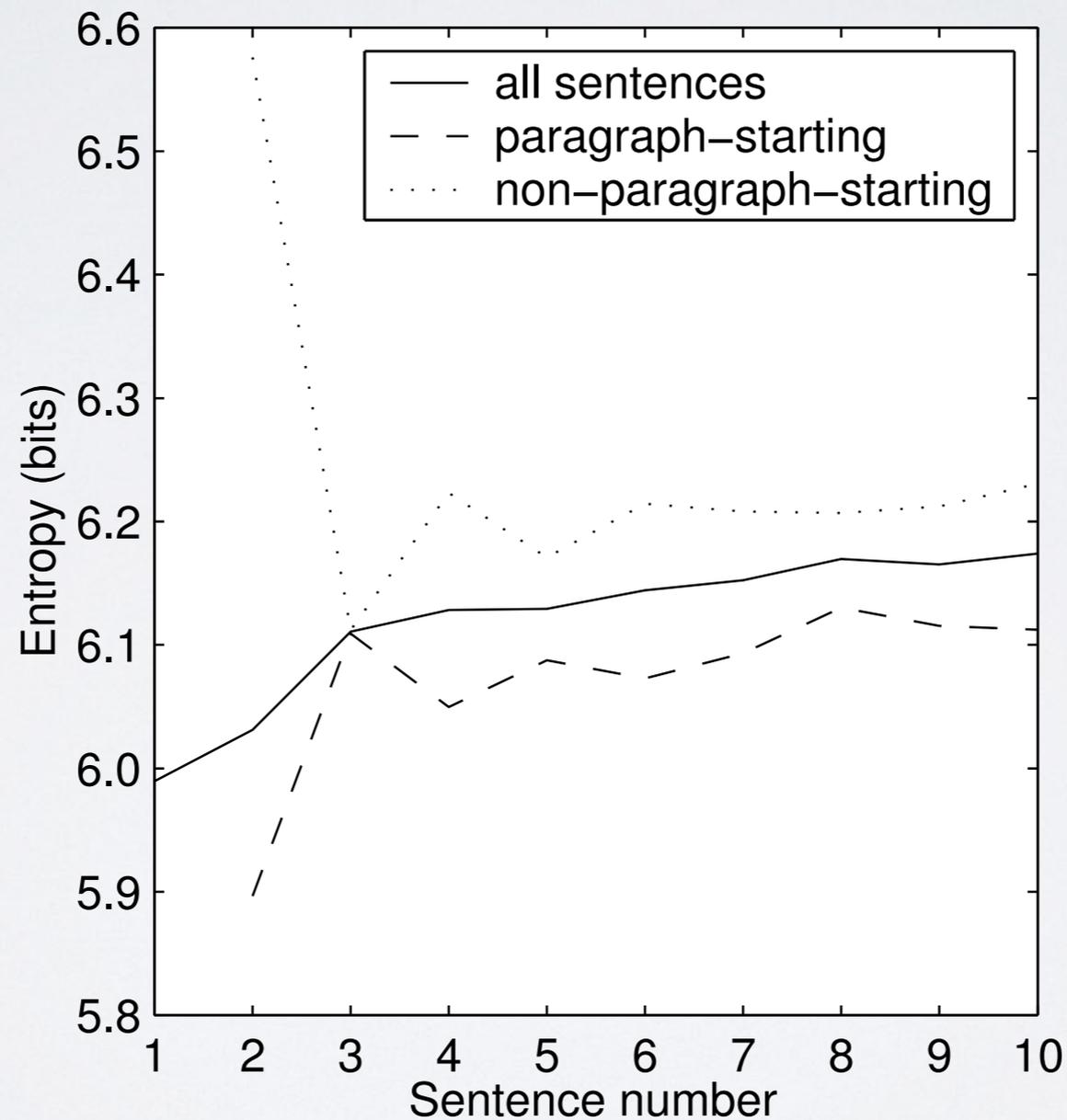


(Jaeger, 2010; Figure 1a)

# Uniform Information Density

## Example 2: Discourse

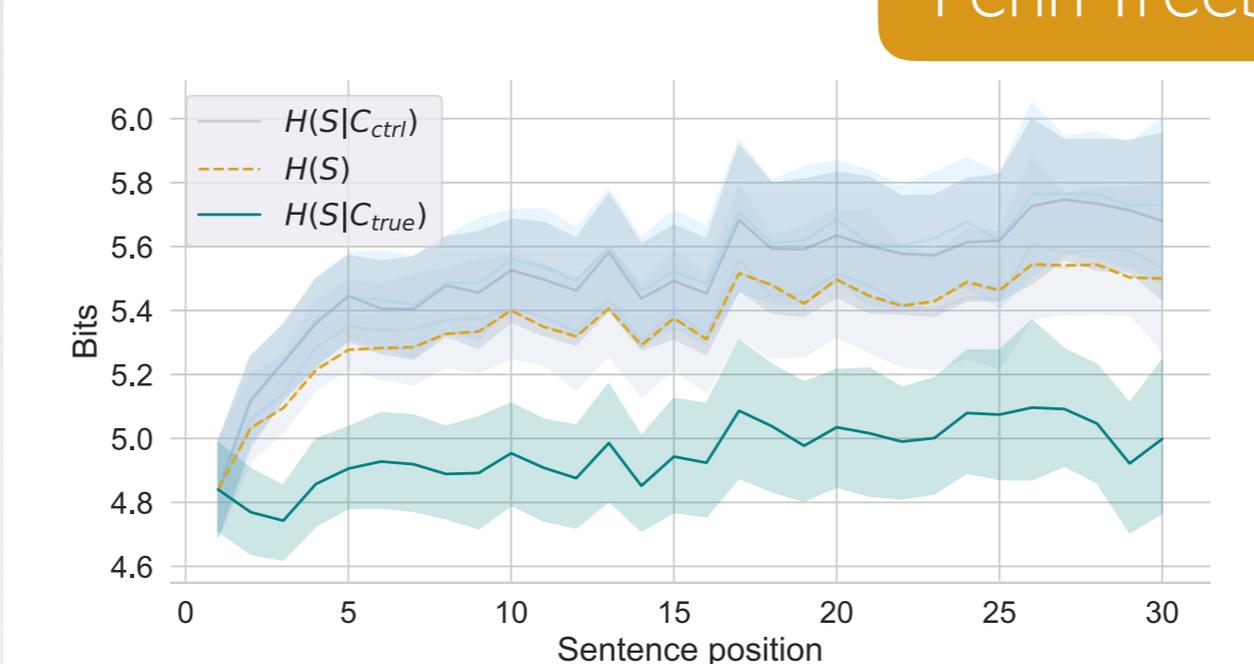
Entropy rate constancy (Genzel & Charniak, 2002, 2003)



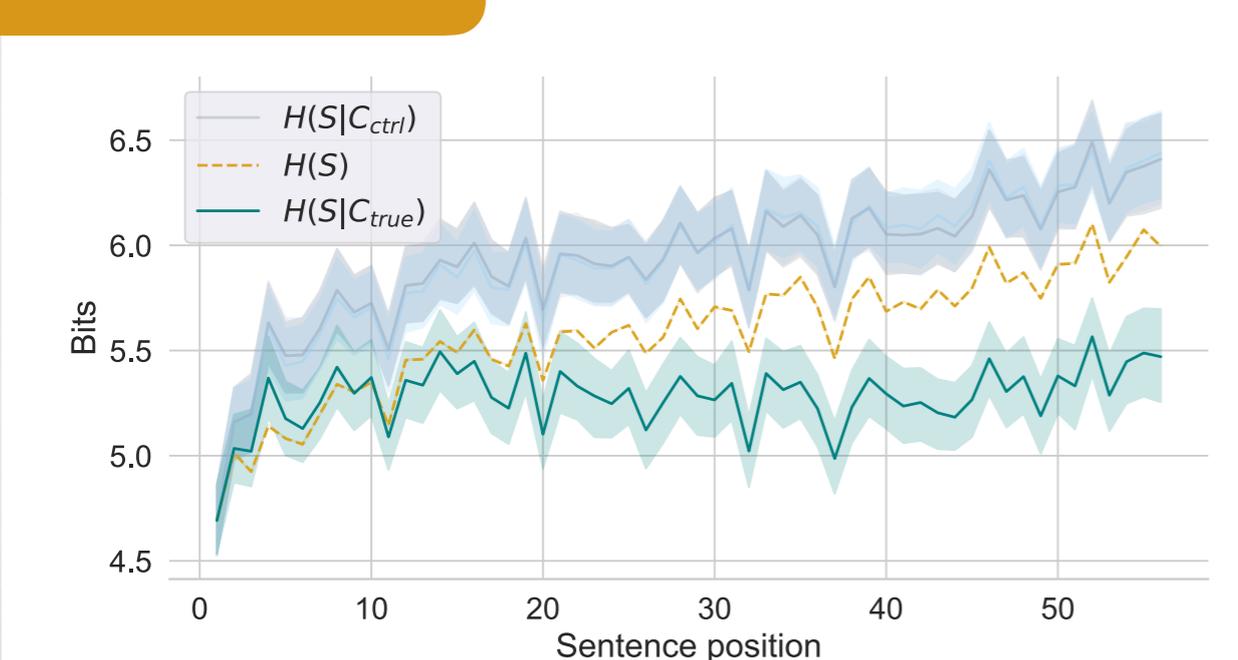
(Genzel & Charniak, 2003; Figure 1)

# Control runs

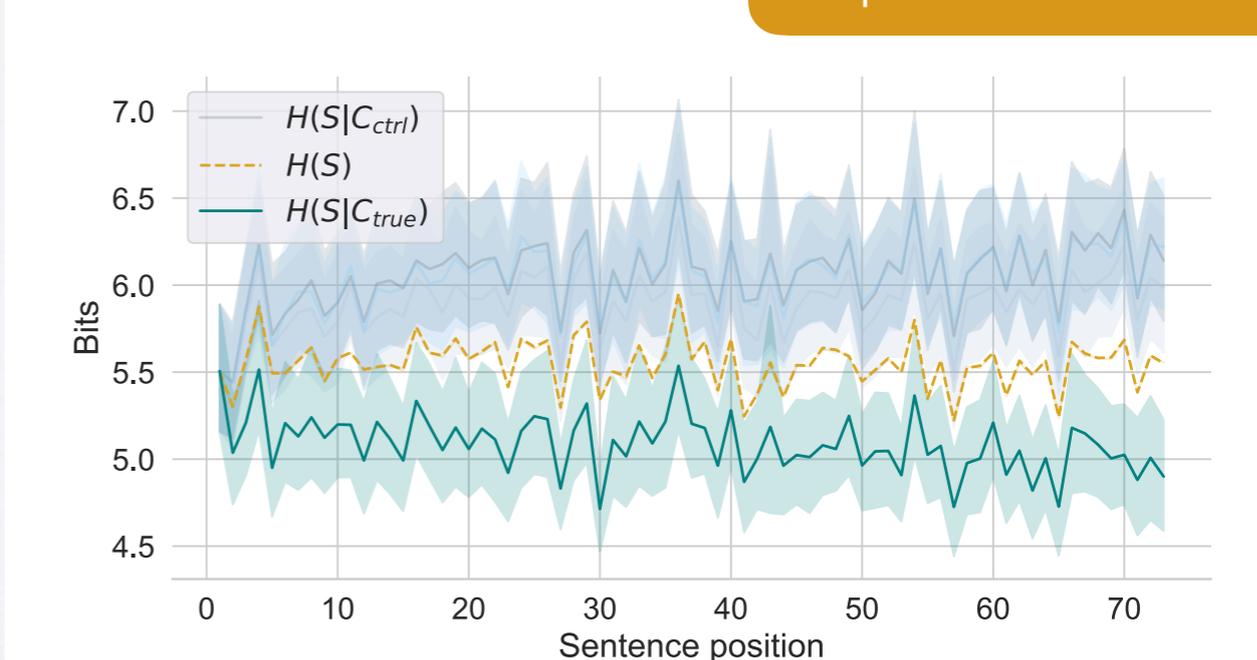
Penn Treebank



PhotoBook

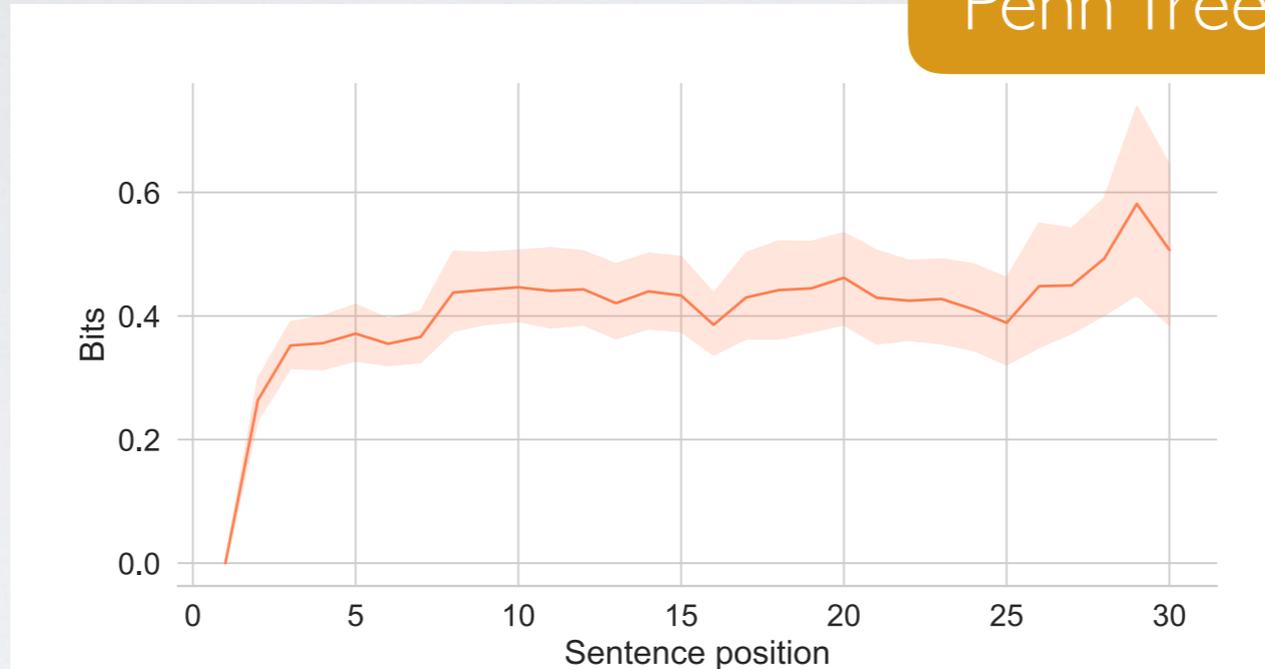


Spoken BNC

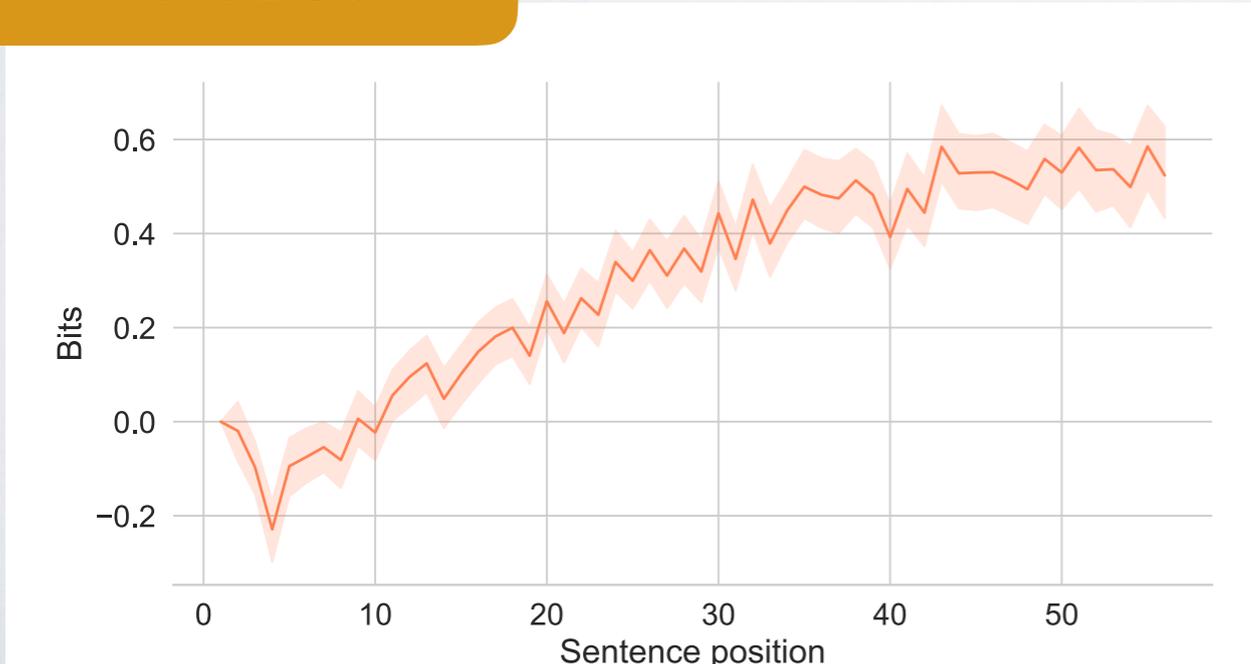


# Context informativeness

Penn Treebank



PhotoBook



Spoken BNC

