

Cross-Linguistic Word Orders Enable an Efficient Tradeoff of Memory and Surprisal

Michael Hahn

Stanford

Judith Degen

Stanford

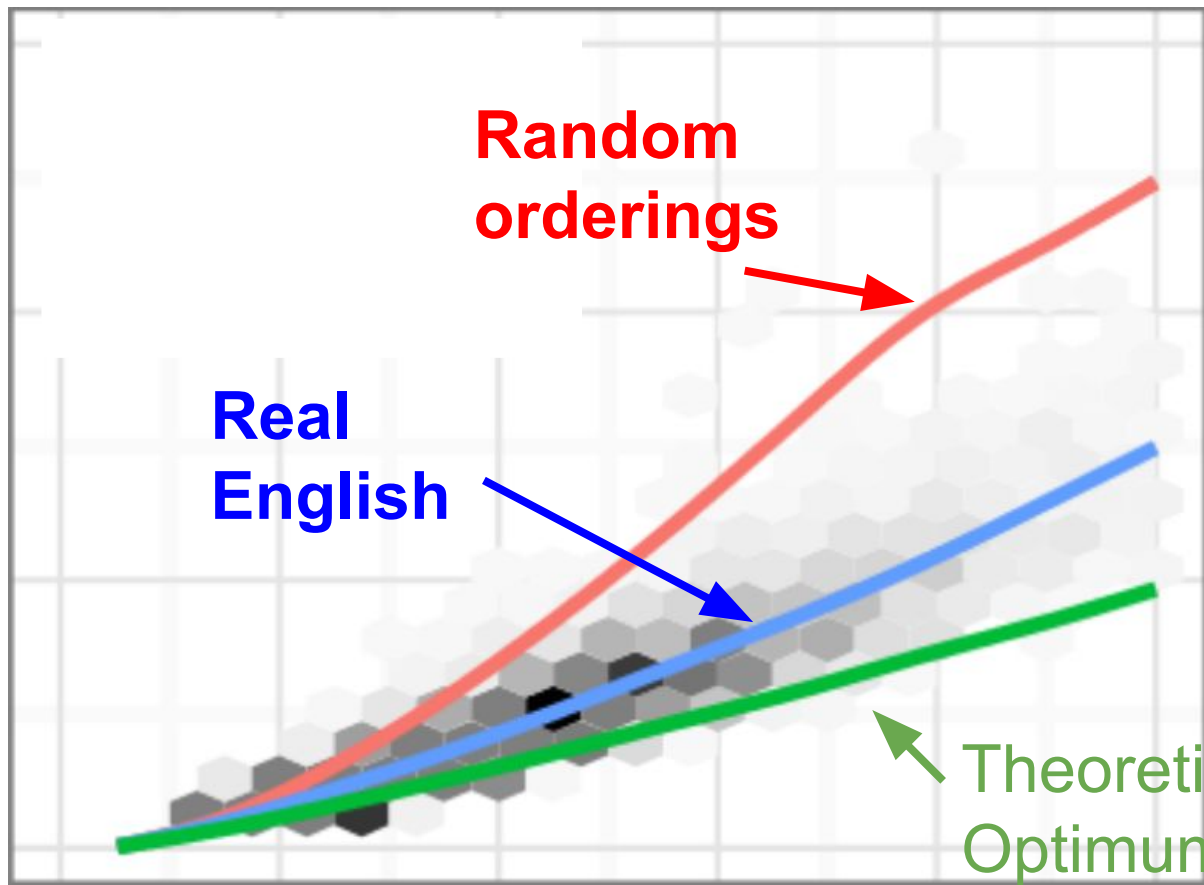
Richard Futrell

UC Irvine

Memory and Word Order

- **Online memory limitations** well-established as a **factor in sentence processing**
- argued to account for **crosslinguistic word order regularities** (Hawkins 1993, Temperley, 2018, ...)

Dependency Length



Random orderings

Real English

Theoretical Optimum

Dependency Length Minimization:

Dependencies are shorter than expected at random

Idea: In certain models, short dependencies reduce memory load (Gibson 1998)

Sentence Length

(Futrell et al., 2015)

Memory and Word Order

- **Online memory limitations** well-established as a **factor in sentence processing**
- argued to account for **crosslinguistic word order regularities** (Hawkins 1993, Temperley, 2018, ...)
- Memory limitations have been cashed out in many different ways
 - Dependency Locality (Gibson 1998)
 - Cue-based retrieval (McElree 2000; Lewis and Vasishth 2005; ...)
 - ...

Memory and Word Order

- **Online memory limitations** well-established as a factor in sentence processing
- argued to account for crosslinguistic word order regularities (Hawkins 1993, Temperley, 2018, ...)
- Memory limitations have been cashed out in many different ways
 - Dependency Locality (Gibson 1998)
 - Cue-based retrieval (McElree 2000; Lewis and Vasishth 2005; ...)
 - ...

Challenge: When testing memory-based explanations of word order, how can we minimize dependence on specific architectural assumptions?

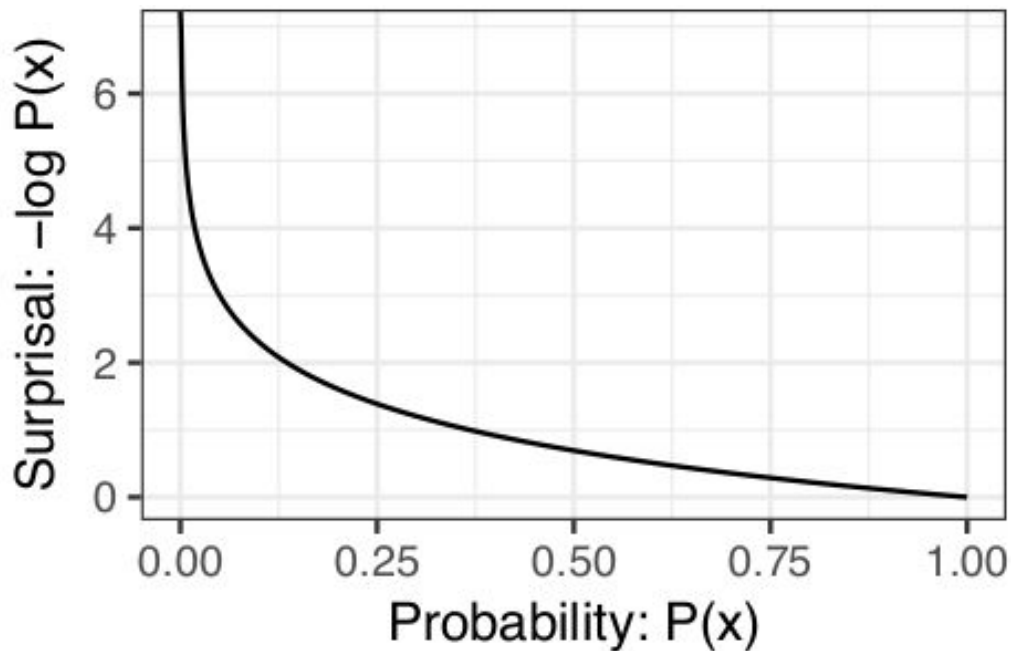
This talk

1. **Information-theoretic formalization** of memory limitations
2. Prove **theorem** describing **tradeoff between memory and surprisal**, without assumptions about memory architecture
3. **Test:** Are crosslinguistic word orders **optimized** for the memory-surprisal tradeoff?

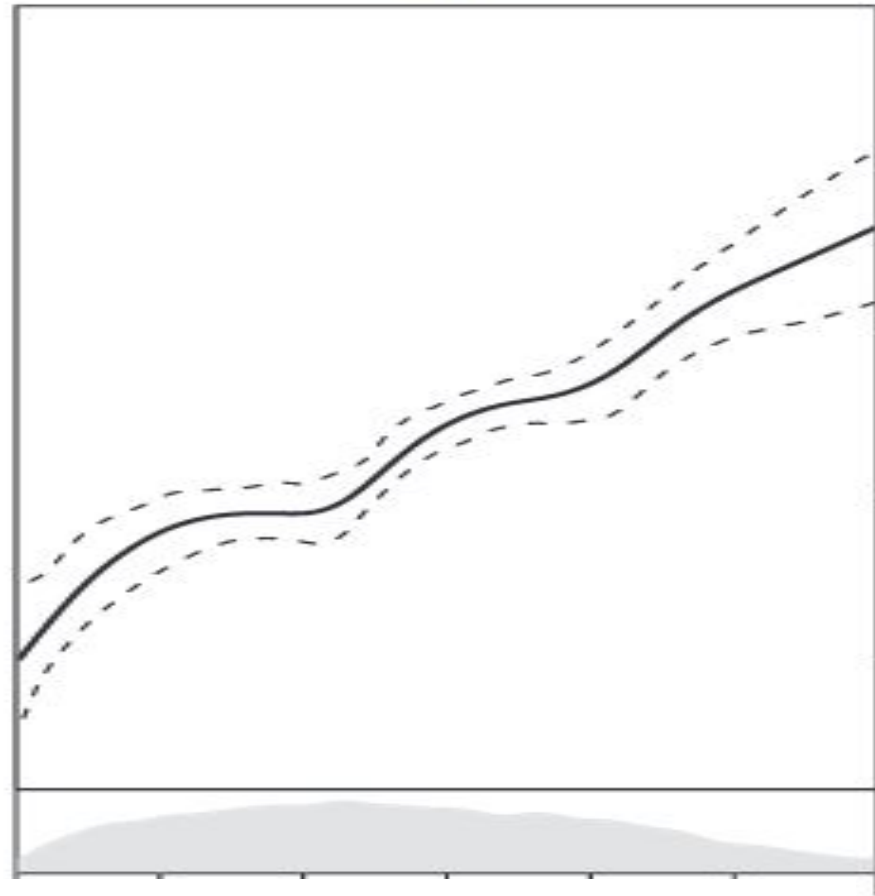
Starting Point: Surprisal Theory (Hale, 2001; Levy, 2008; Smith & Levy, 2013; Hale, 2016)

Processing difficulty at a word is equal to the surprisal of that word in context:

$$C(w \mid \text{context}) = -\log P(w \mid \text{context})$$



Reading Time



Surprisal

(Smith and Levy 2013)

Surprisal

Hey! What's



up?	0.65
the	0.2
a	0.15
...	...



Surprisal

Hey! What's up?



up?	0.65
the	0.2
a	0.15
...	...

Surprisal(up|Hey! What's)

$$= -\log 0.65 \sim 0.18$$



Surprisal

Hey! What's



Hey! What's	up?	0.65
	the	0.2
	a	0.15

Surprisal(up|Hey! What's)

$$= -\log 0.65 \sim 0.18$$



Surprisal

Listener has forgotten the past.

Hey! What's

???? ???? →

the	0.23
a	0.2
in	0.15
up	0.09



Surprisal(up|???)

= $-\log 0.09 \sim 2.4$



Surprisal

Hey! What's

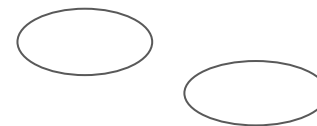


Listener has forgotten the past.

???? ???? |

the	0.23
a	0.2
in	0.15
up	0.09

Cannot utilize context for prediction.



Surprisal(up|???)

= $-\log 0.09 \sim 2.4$

Surprisal

Listener has forgotten the past.

Cannot utilize context for prediction.

Hey! What's

???? ????	the	0.23
	a	0.2
	in	0.15
	up	0.09

Incurs higher surprisal

$$\text{Surprisal}(\text{up}|\text{???)}$$

$$= -\log 0.09 \sim 2.4$$



Surprisal

Hey! What's



???? ????	the	0.23
	a	0.2
	in	0.15
	up	0.09

A forgetful listener incurs higher average surprisal.

Surprisal(up|???)

= $-\log 0.09 \sim 2.4$

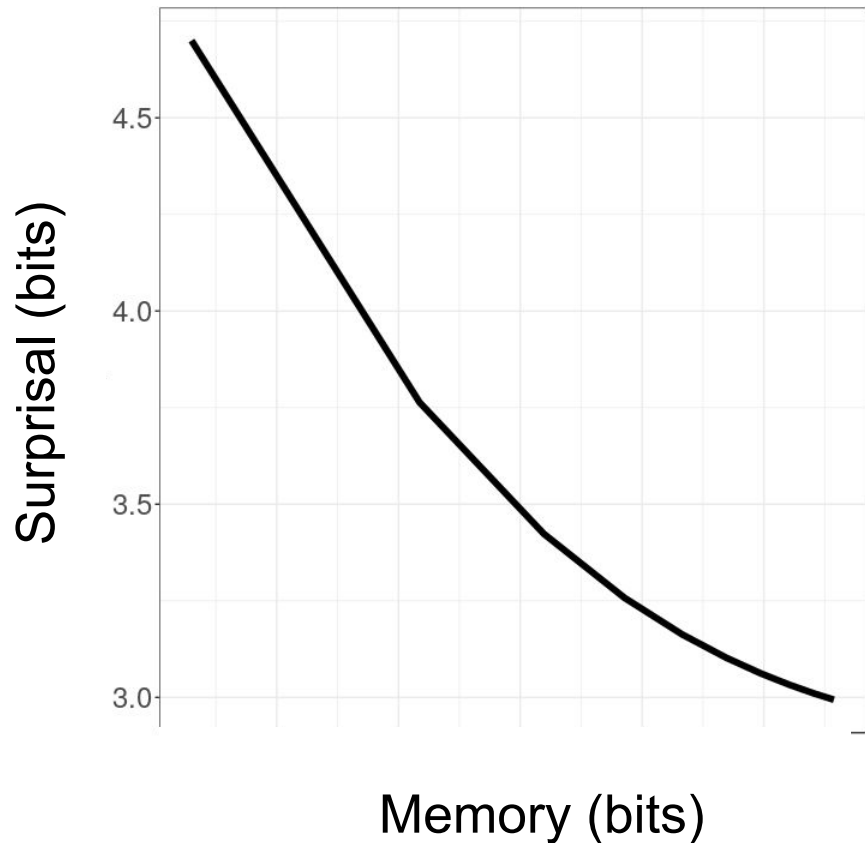


Memory-Surprisal Tradeoff

Having better
representation of the
past improves
prediction of the future
on average.

Memory-Surprisal Tradeoff

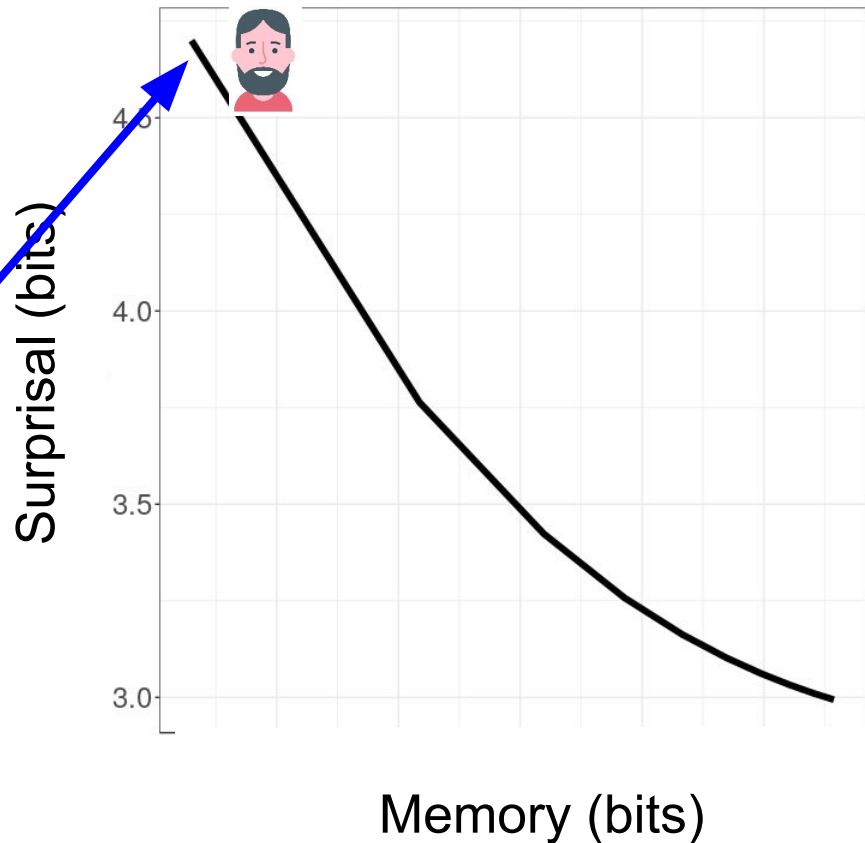
Having better representation of the past improves prediction of the future on average.



Memory-Surprisal Tradeoff

Having better representation of the past improves prediction of the future on average.

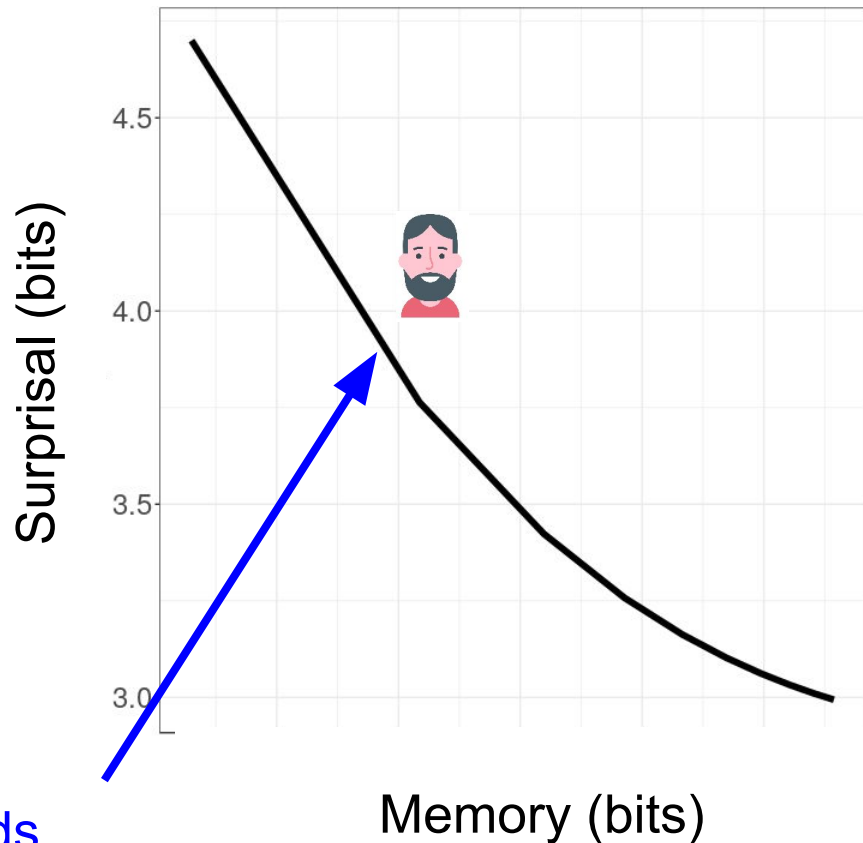
Remembering 0 bits leads to maximum surprisal



Memory-Surprisal Tradeoff

Having better representation of the past improves prediction of the future on average.

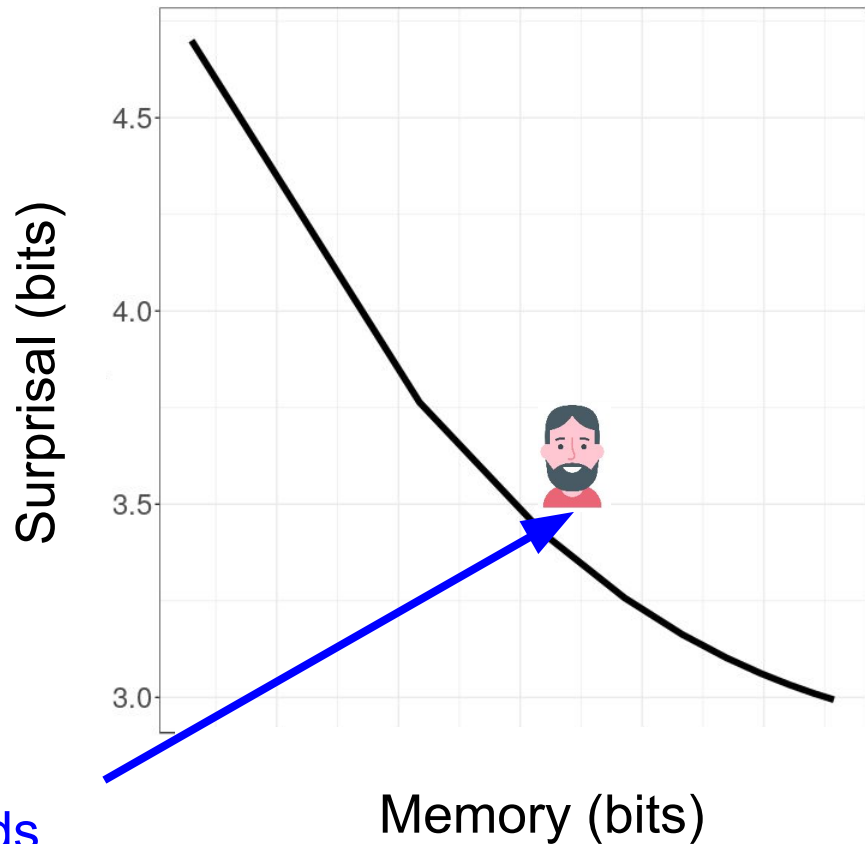
Remembering more leads to lower surprisal



Memory-Surprisal Tradeoff

Having better representation of the past improves prediction of the future on average.

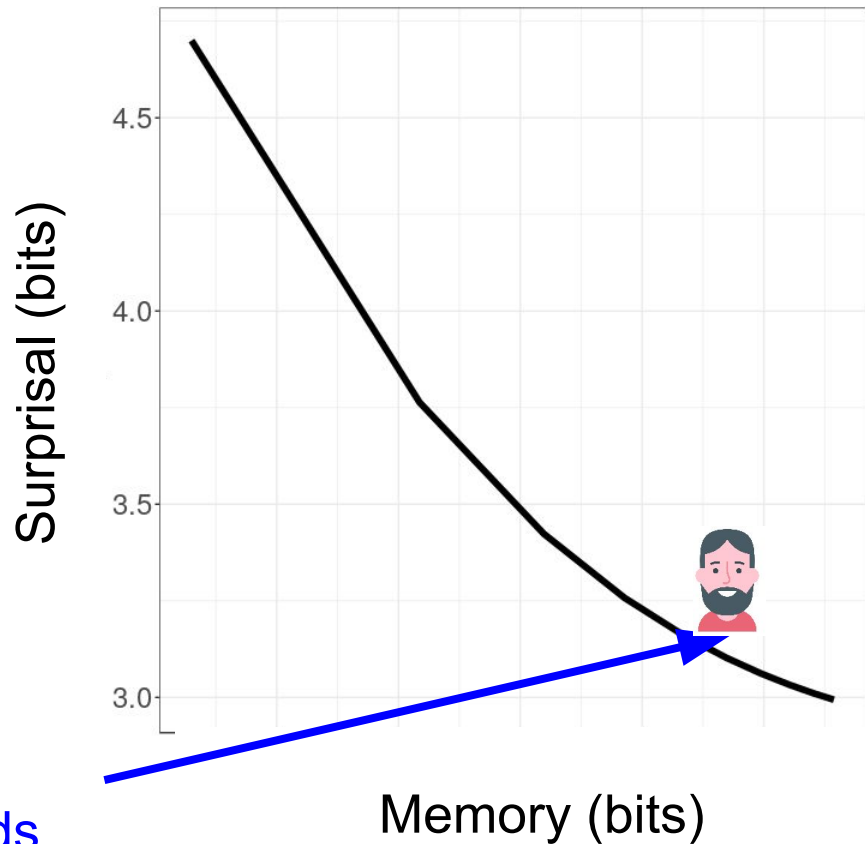
Remembering more leads to lower surprisal



Memory-Surprisal Tradeoff

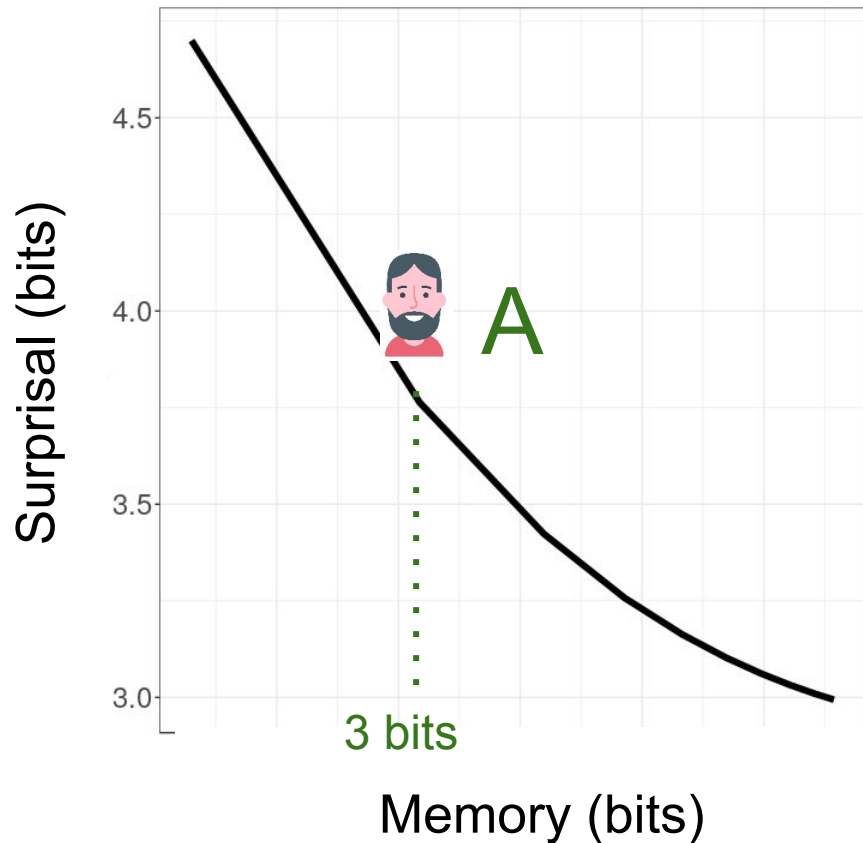
Having better representation of the past improves prediction of the future on average.

Remembering more leads to lower surprisal



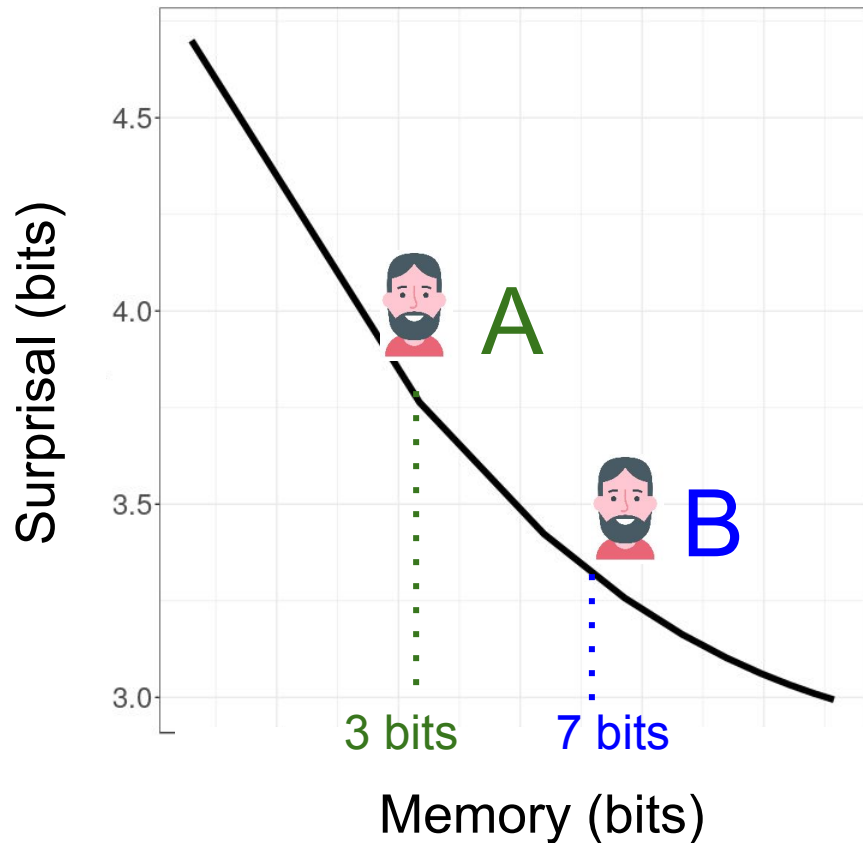
Memory-Surprisal Tradeoff

Having better representation of the past improves prediction of the future on average.



Memory-Surprisal Tradeoff

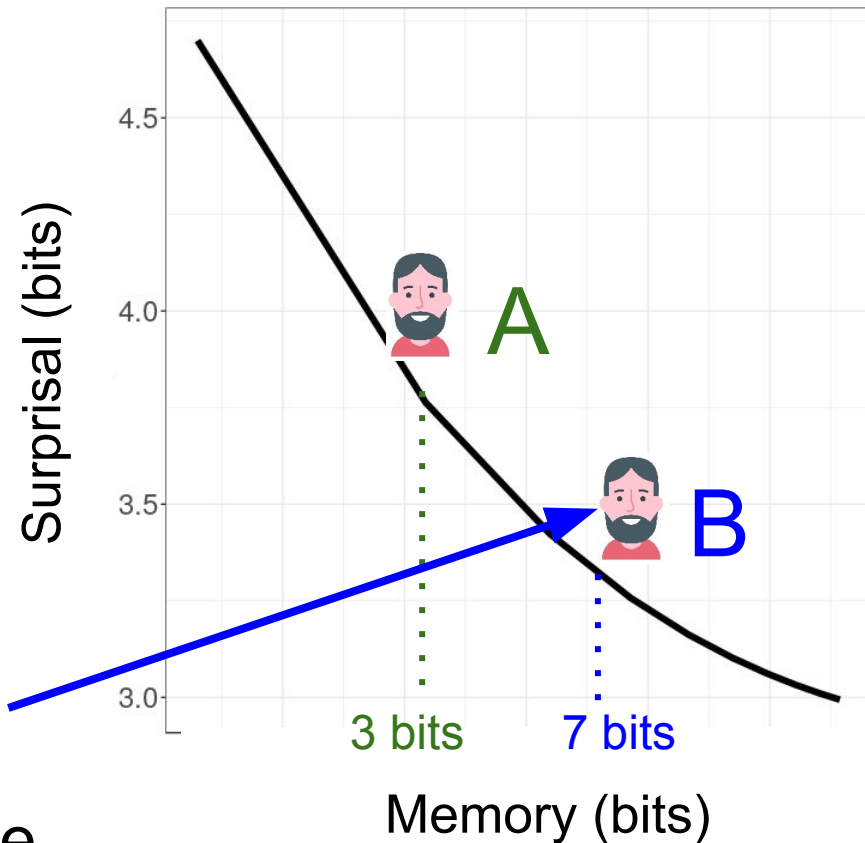
Having better representation of the past improves prediction of the future on average.



Memory-Surprisal Tradeoff

Having better representation of the past improves prediction of the future on average.

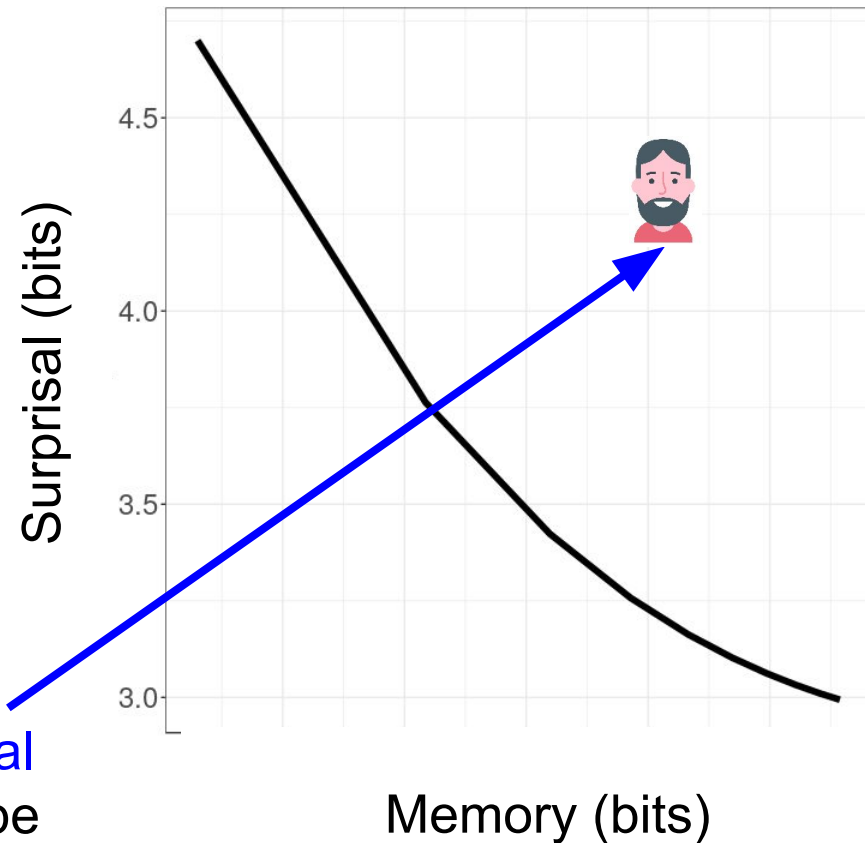
B will incur **lower surprisal** on average



Memory-Surprisal Tradeoff

Having better representation of the past improves prediction of the future on average.

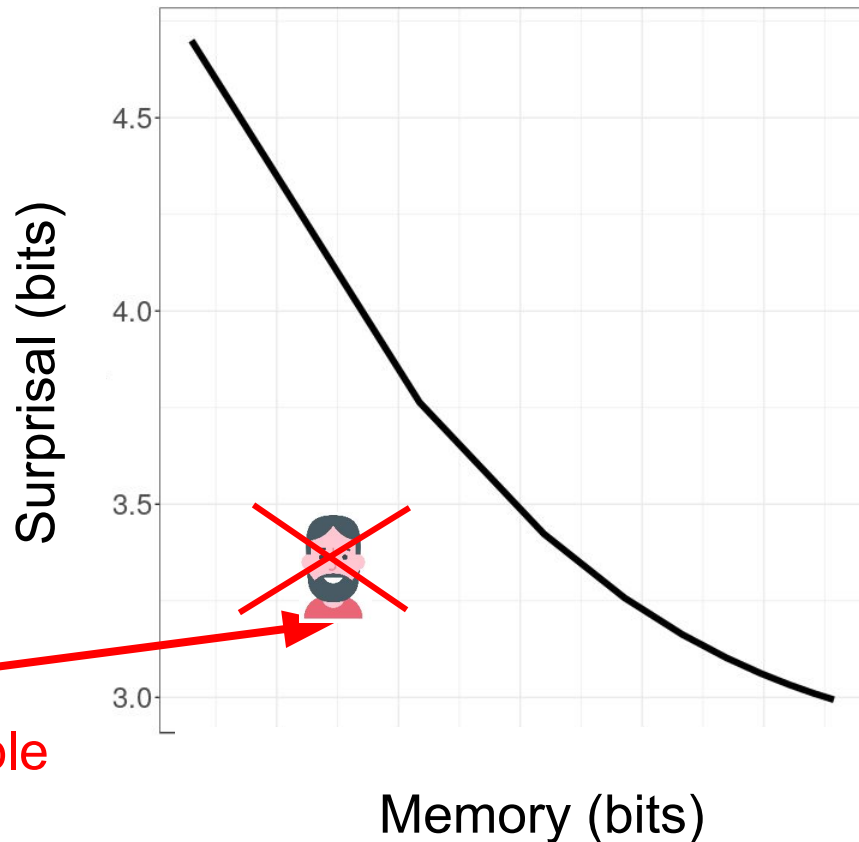
A listener with **suboptimal** memory allocation can be **above** the curve



Memory-Surprisal Tradeoff

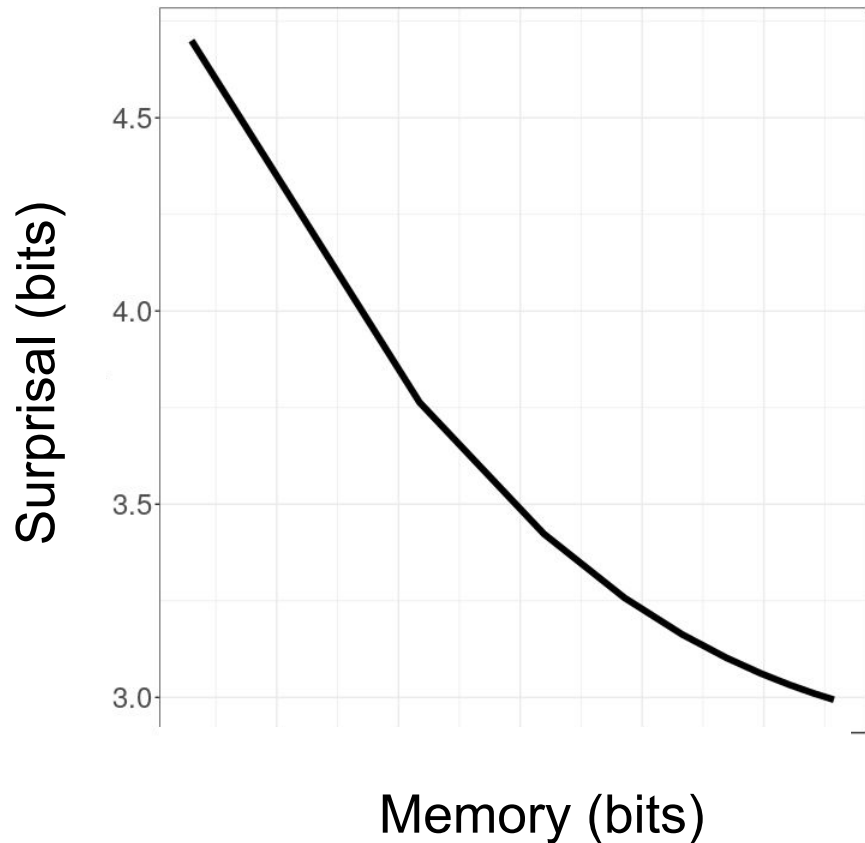
Having better representation of the past improves prediction of the future on average.

Mathematically impossible to be below

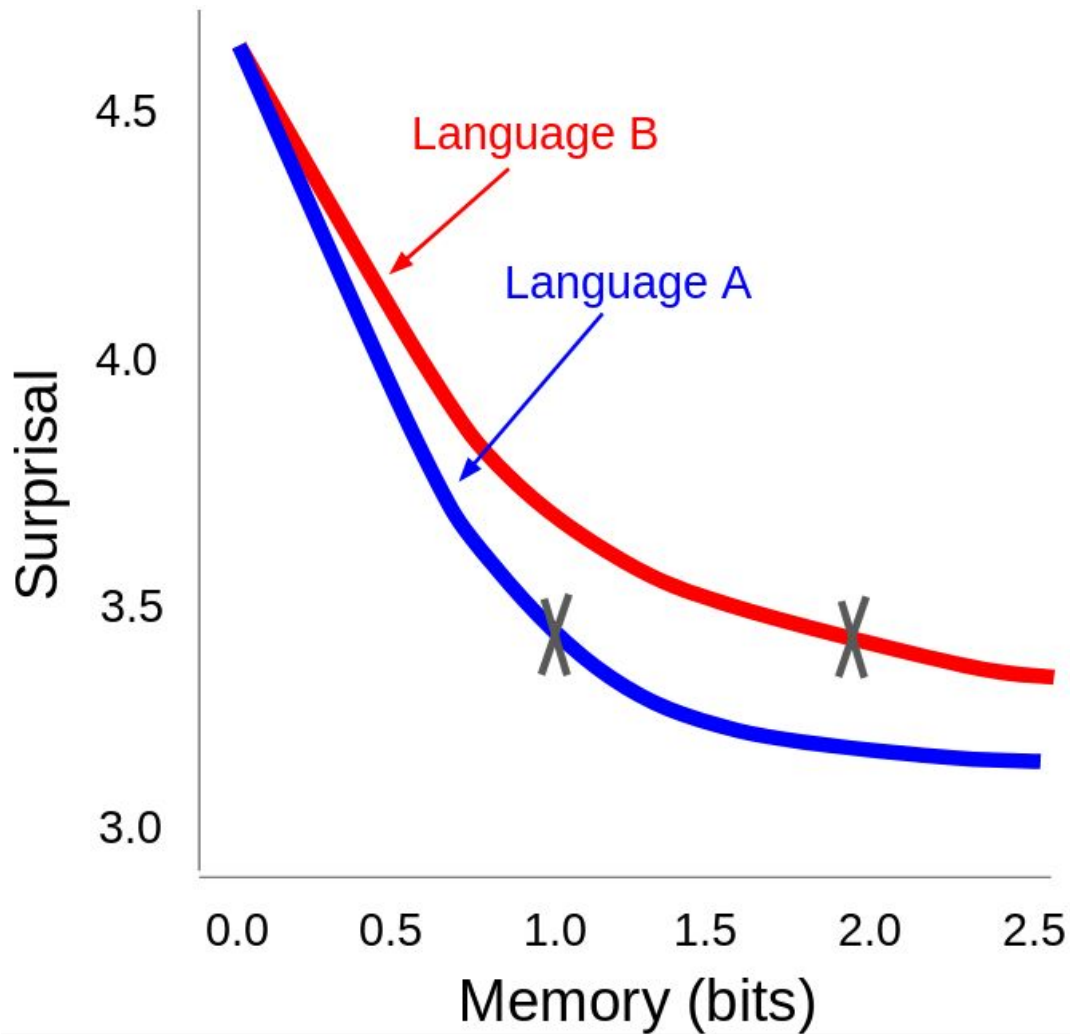


Memory-Surprisal Tradeoff

Having better representation of the past improves prediction of the future on average.

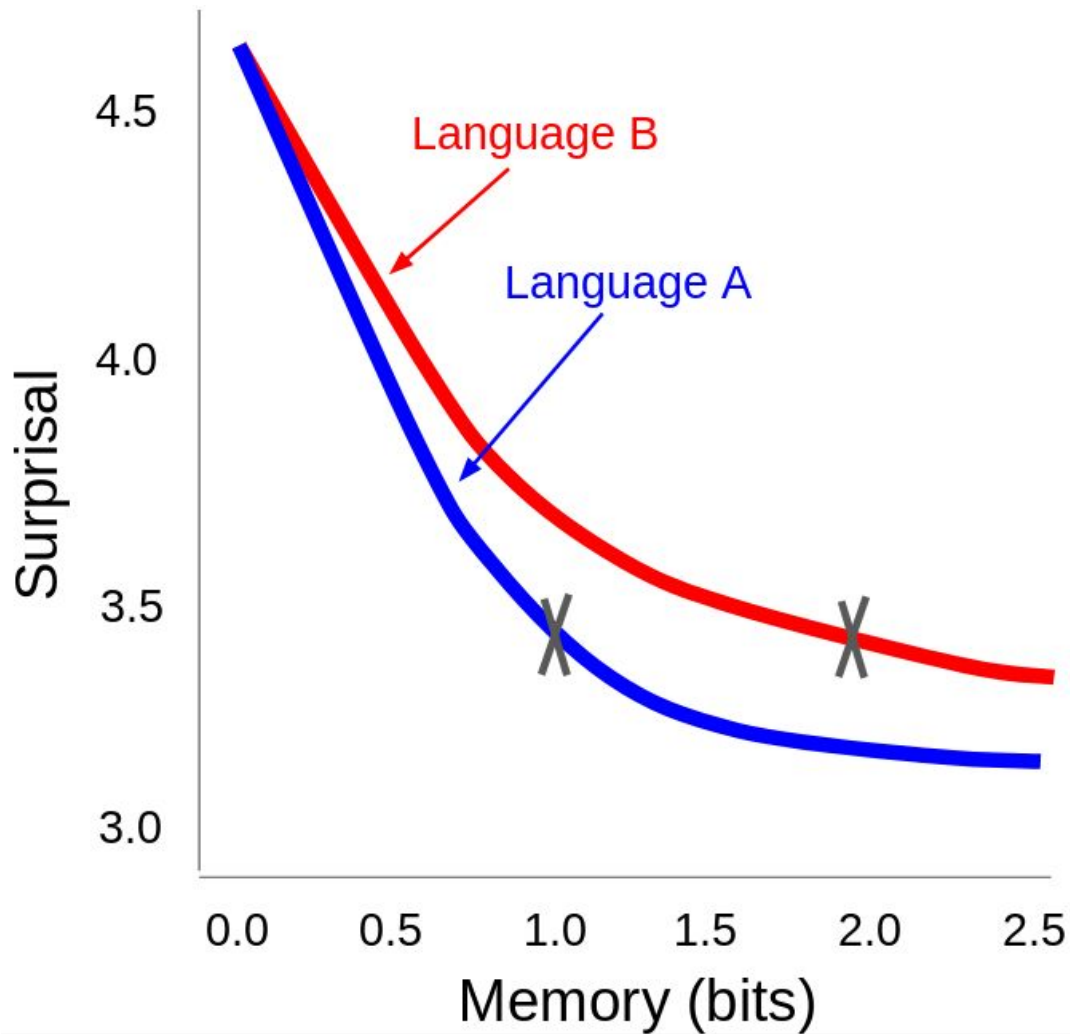


Different languages can lead to different tradeoffs



Different languages can lead to different tradeoffs

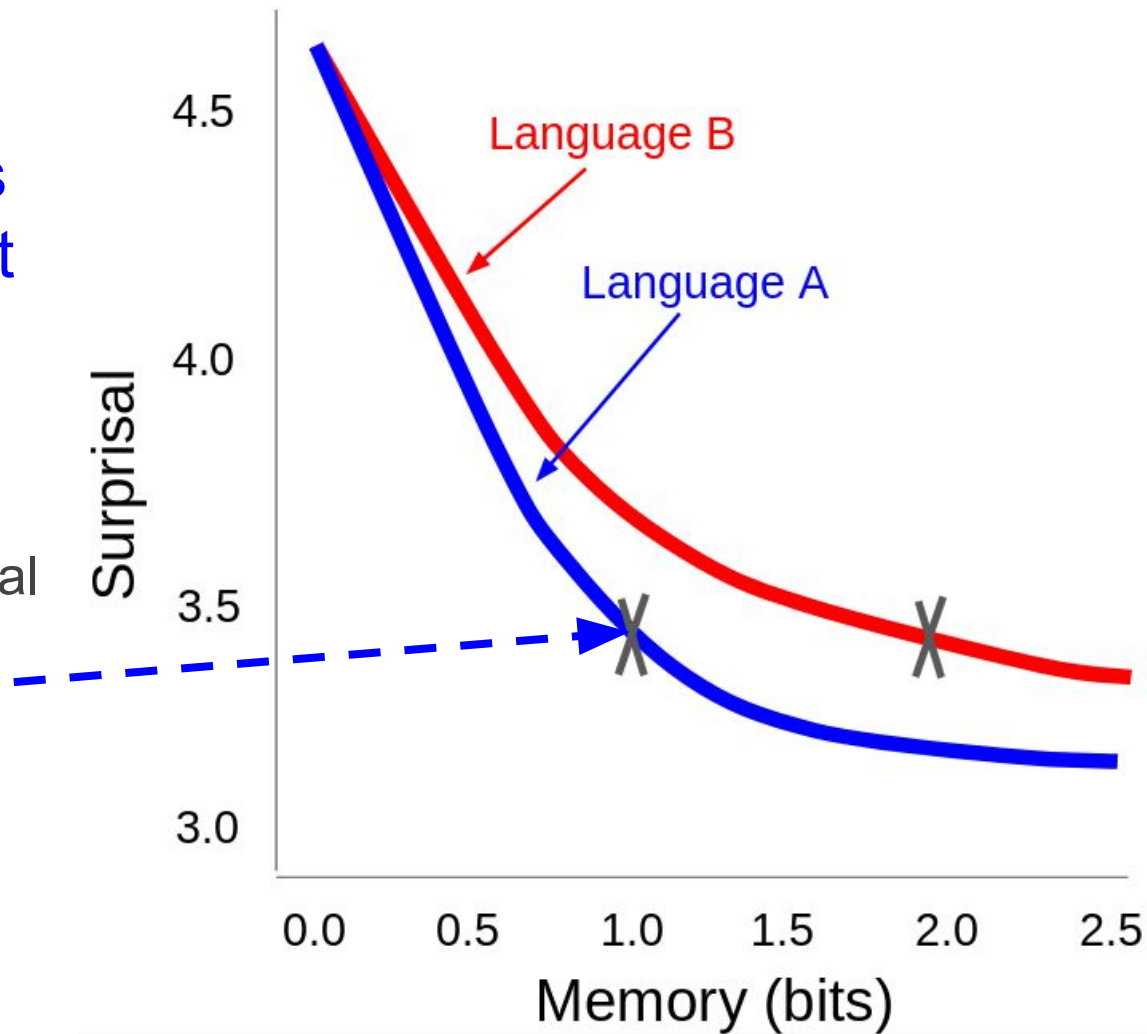
Achieving at most **3.5 bits** of average surprisal takes...



Different languages can lead to different tradeoffs

Achieving at most **3.5 bits** of average surprisal takes...

1.0 bit of memory in Language A

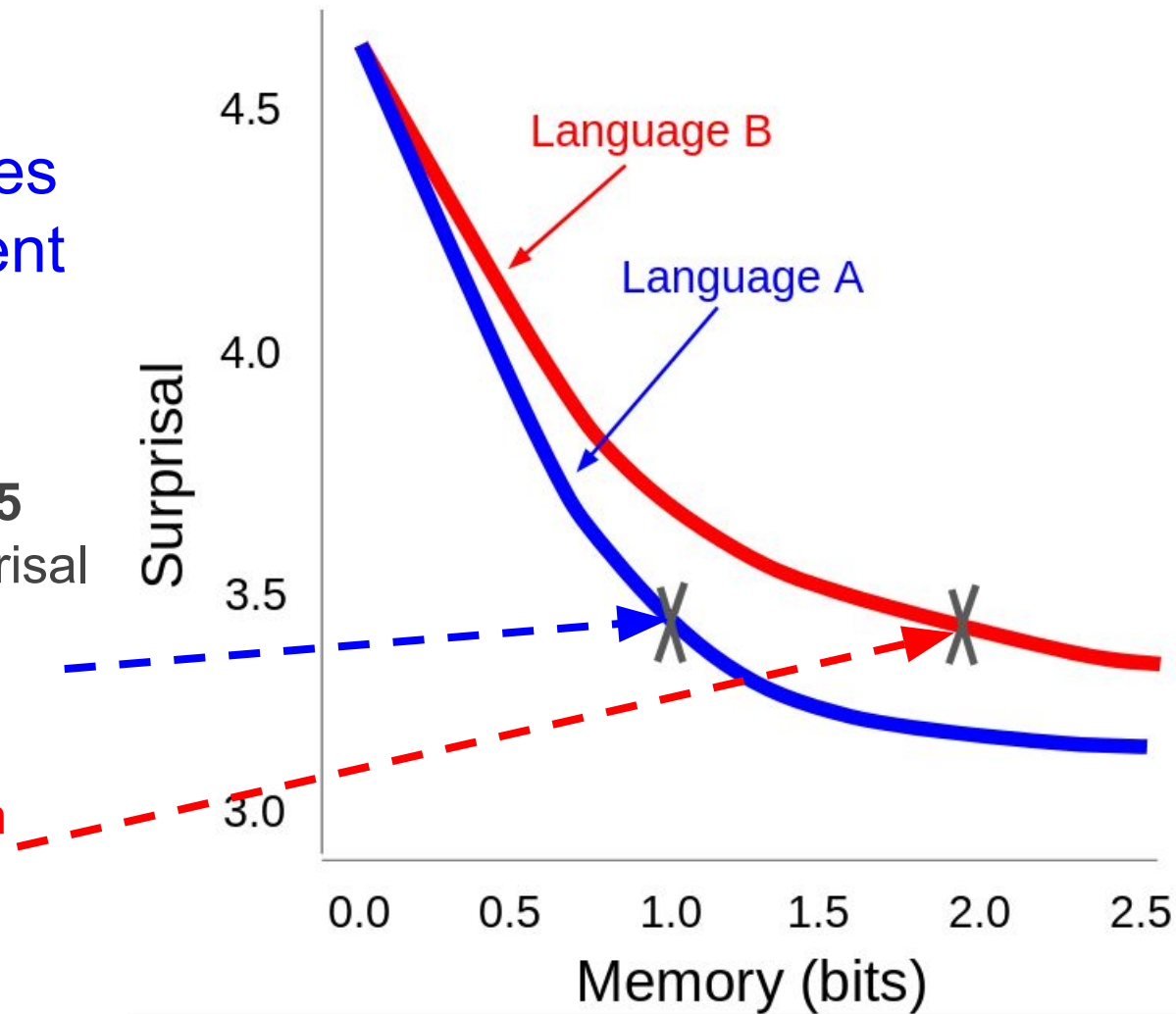


Different languages can lead to different tradeoffs

Achieving at most **3.5 bits** of average surprisal takes...

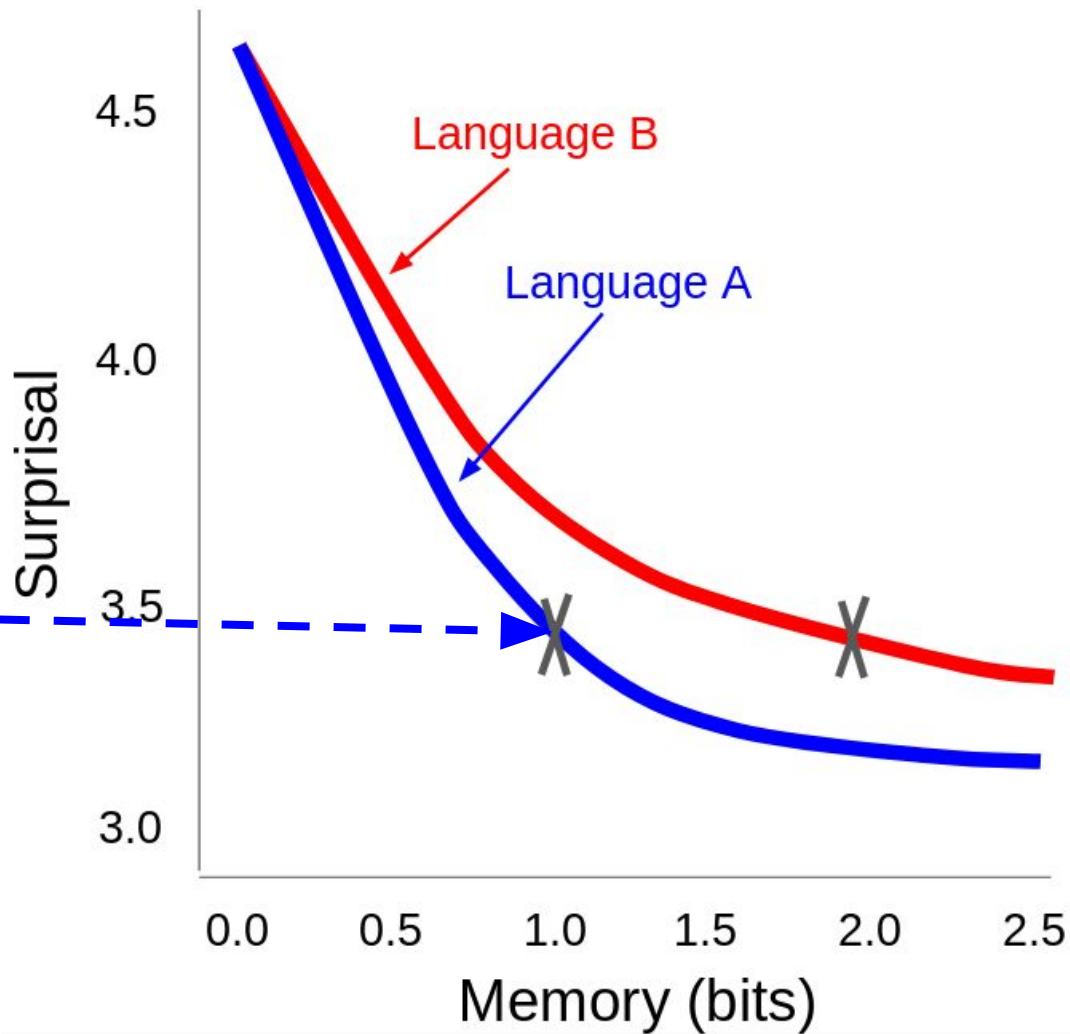
1.0 bit of memory in Language A

2.0 bits of memory in Language B



Different languages can lead to different tradeoffs

Language A has a more favorable tradeoff



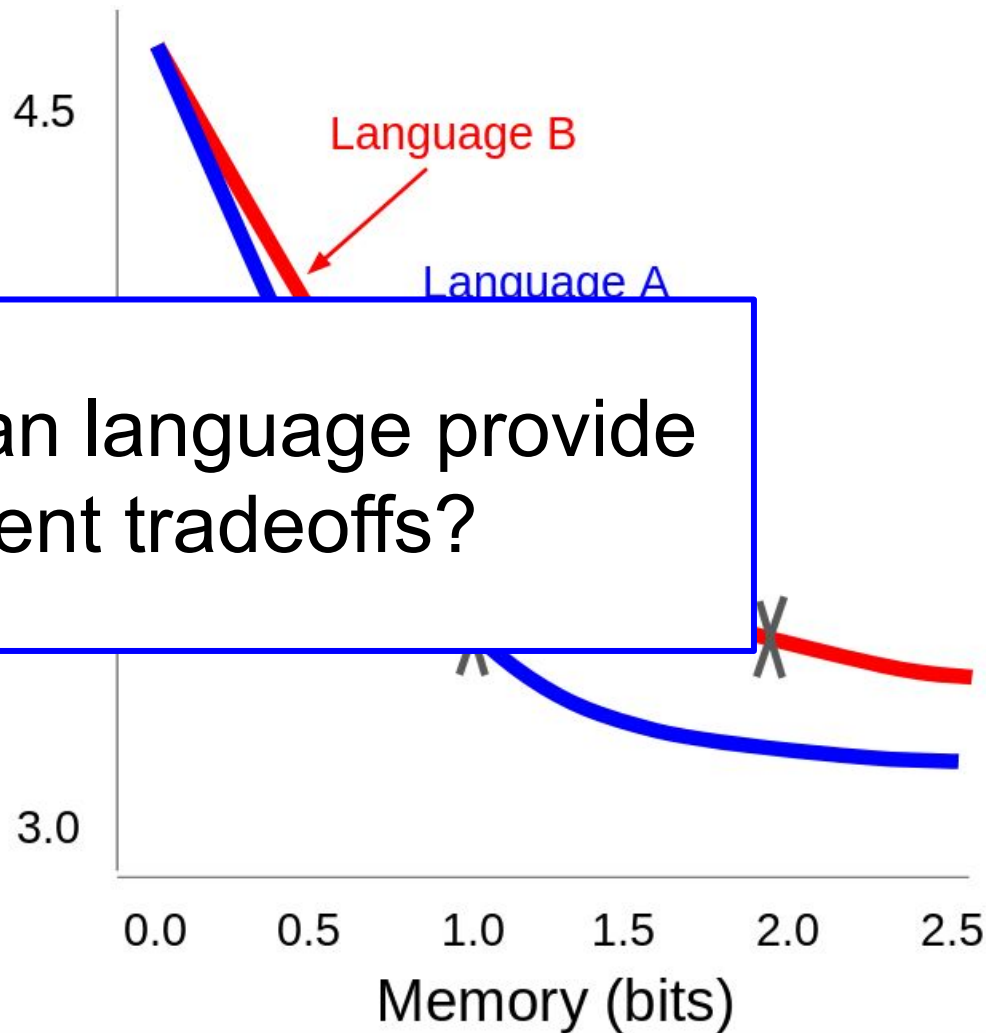
This talk

1. [Information-theoretic formalization](#) of memory limitations
2. Prove theorem describing tradeoff between memory and surprisal, without assumptions about memory architecture
3. Test: Are crosslinguistic word orders optimized for the memory-surprisal tradeoff?

Different languages
can lead to different
tradeoffs

Language
a more fa
tradeoff

Does human language provide
efficient tradeoffs?



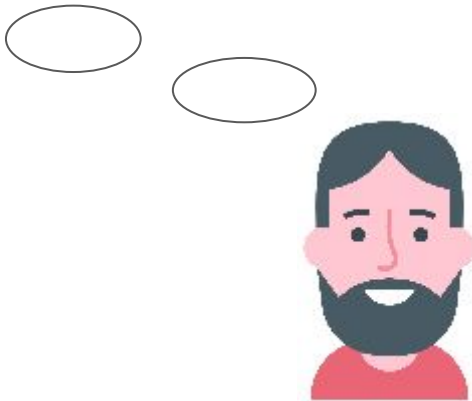
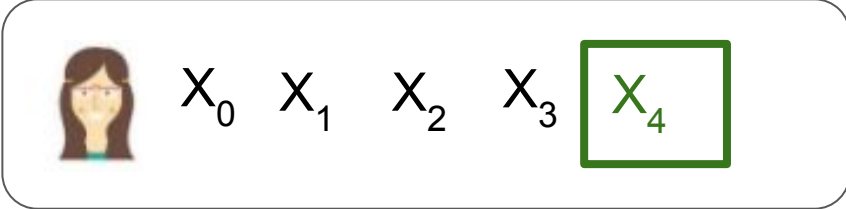
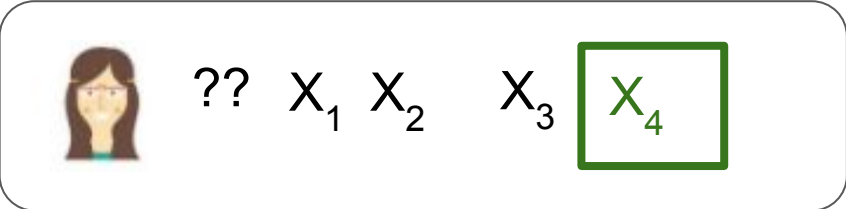
This talk

1. Information-theoretic formalization of memory limitations
2. Prove **theorem** describing **tradeoff between memory and surprisal**, without assumptions about memory architecture
3. Test: Are crosslinguistic word orders optimized for the memory-surprisal tradeoff?

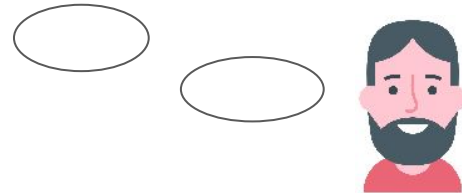
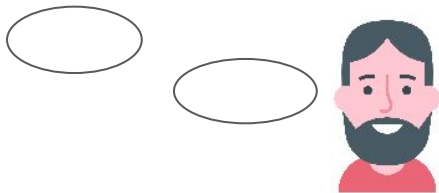
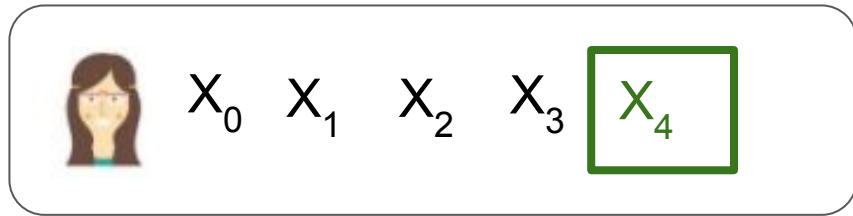
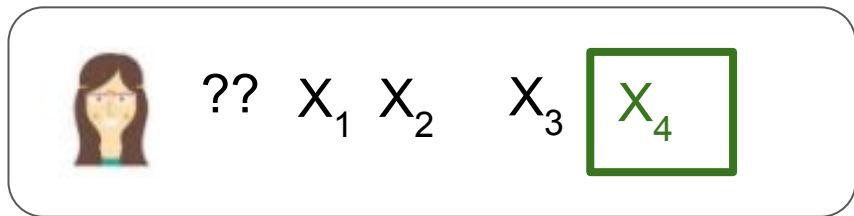
Conditional Mutual Information

$$I[X_t, X_0 | X_1, \dots, X_{t-1}]$$

Conditional Mutual Information



Conditional Mutual Information



Surprisal based on t words of context

MINUS

Surprisal based on $t+1$ words of context

is the definition of:

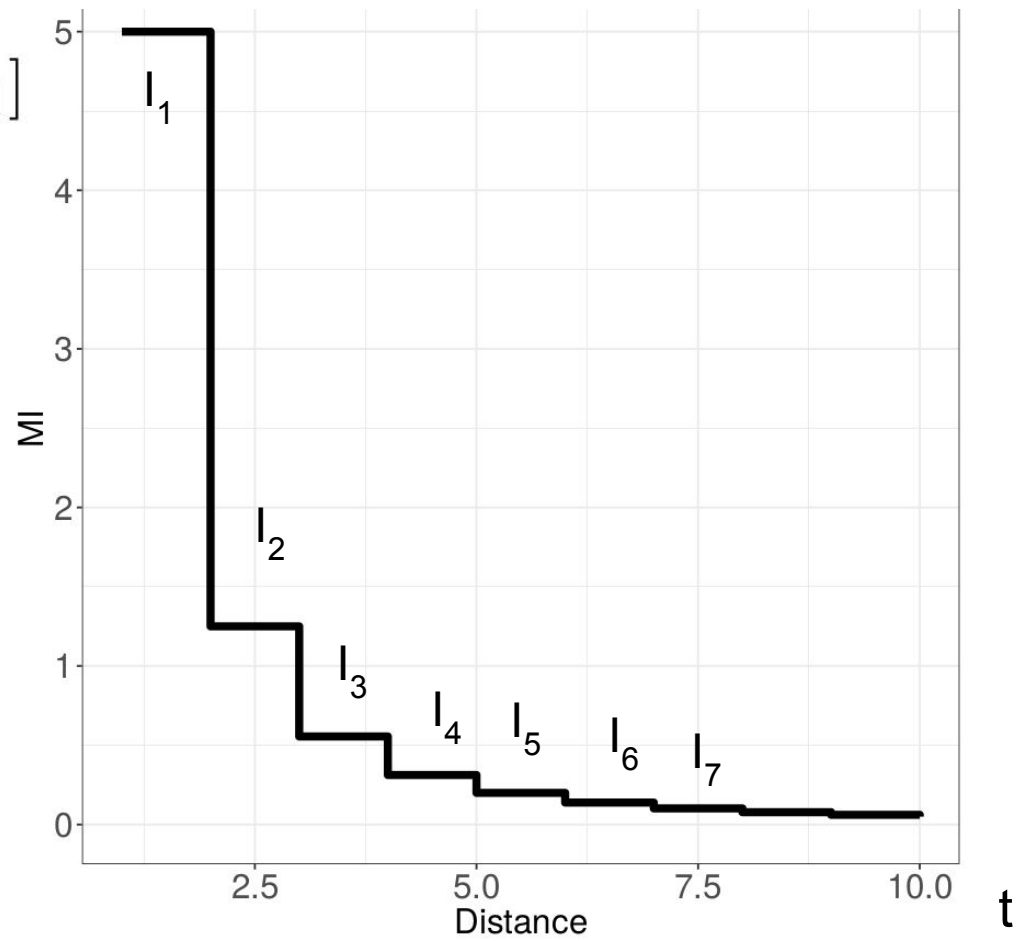
$$I[X_t, X_0 | X_1, \dots, X_{t-1}]$$

Conditional Mutual Information

How much information do words t steps apart contain about each other, **controlling for info redundant with intervening words?**

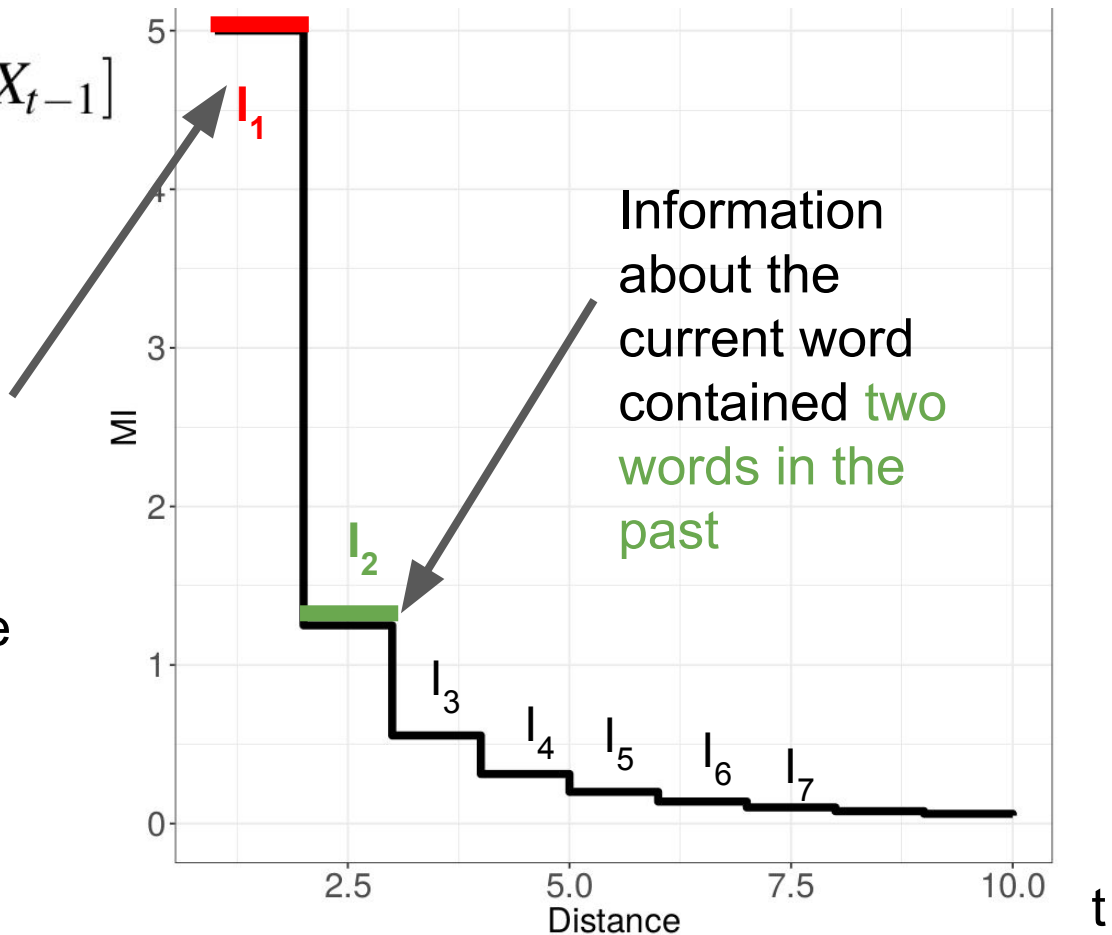
$$I[X_t, X_0 | X_1, \dots, X_{t-1}]$$

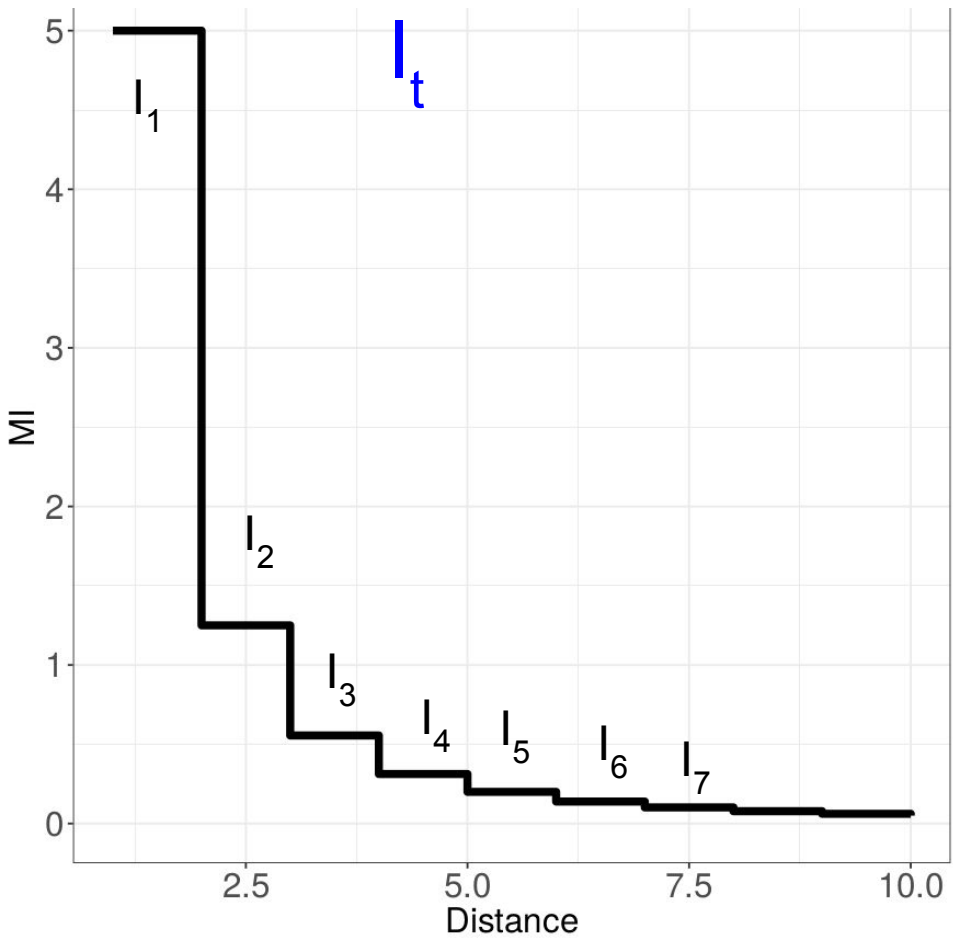
$$I[X_t, X_0 | X_1, \dots, X_{t-1}]$$

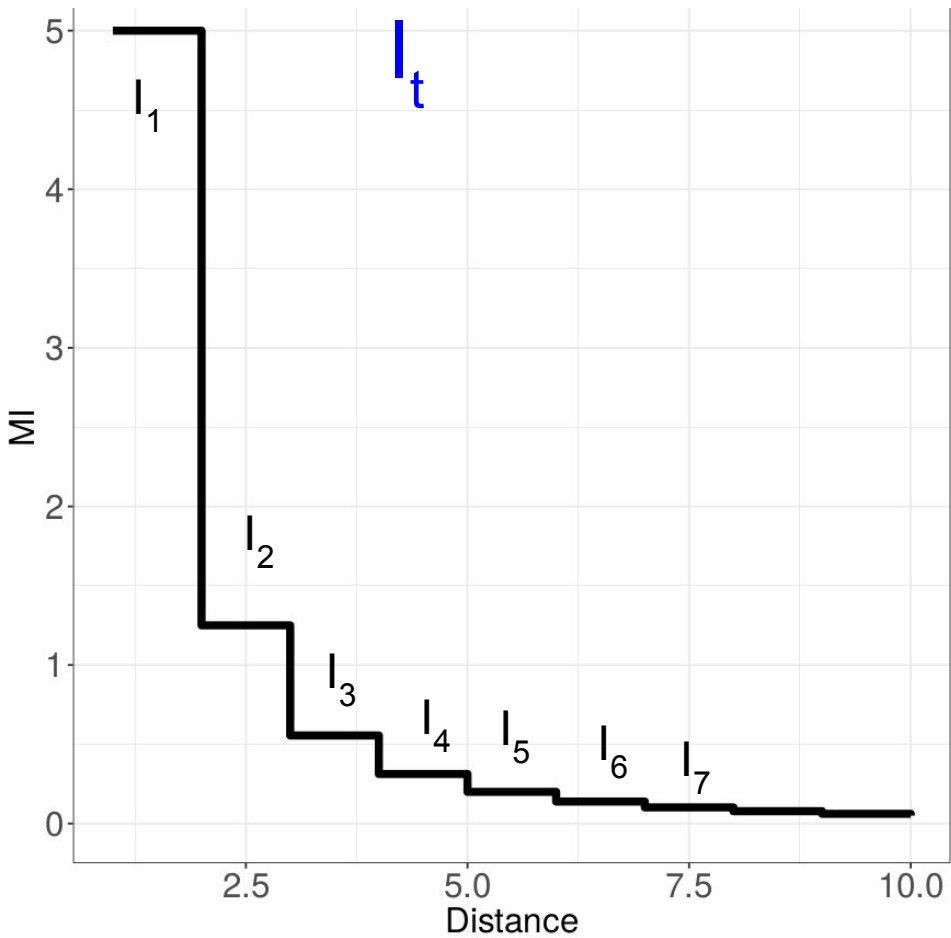


$$I[X_t, X_0 | X_1, \dots, X_{t-1}]$$

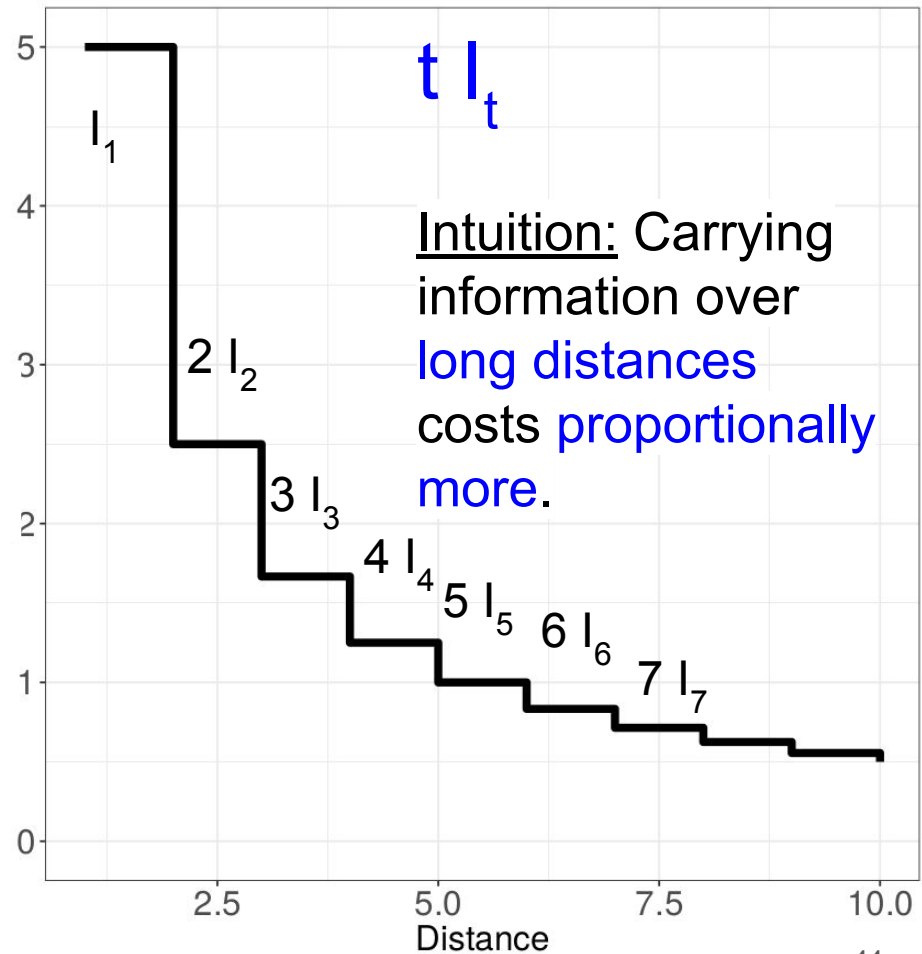
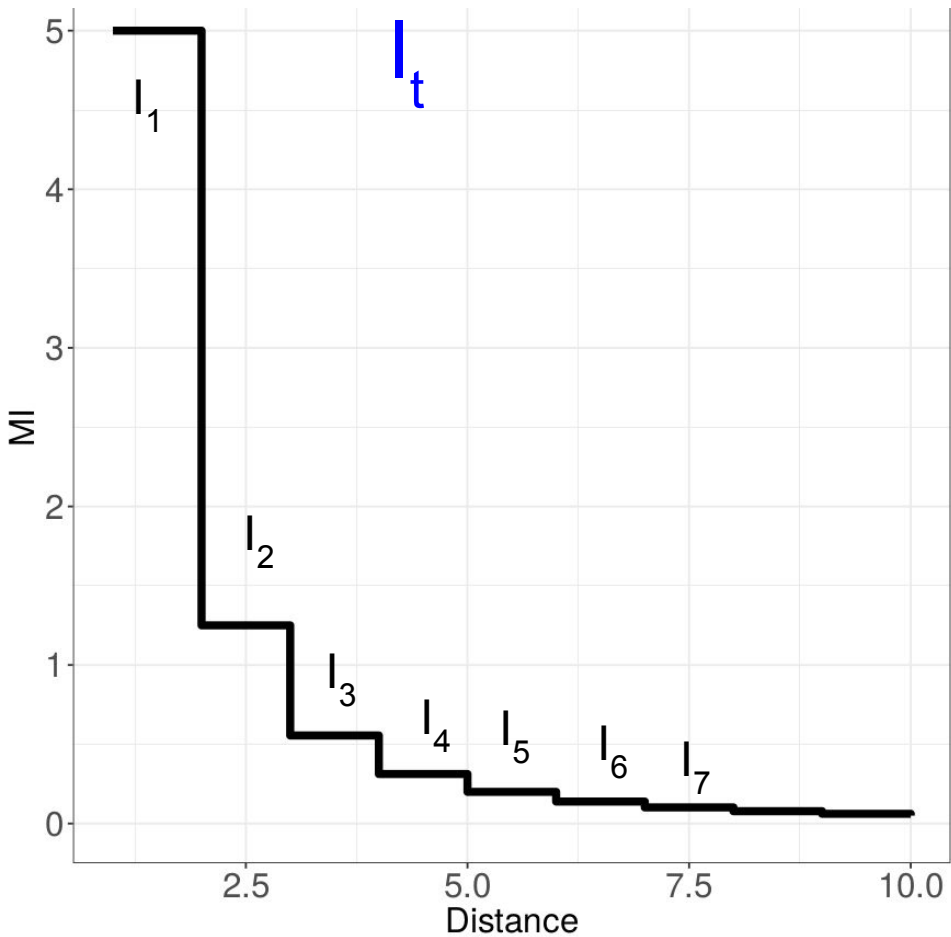
Information about the current word contained in the **last preceding word**





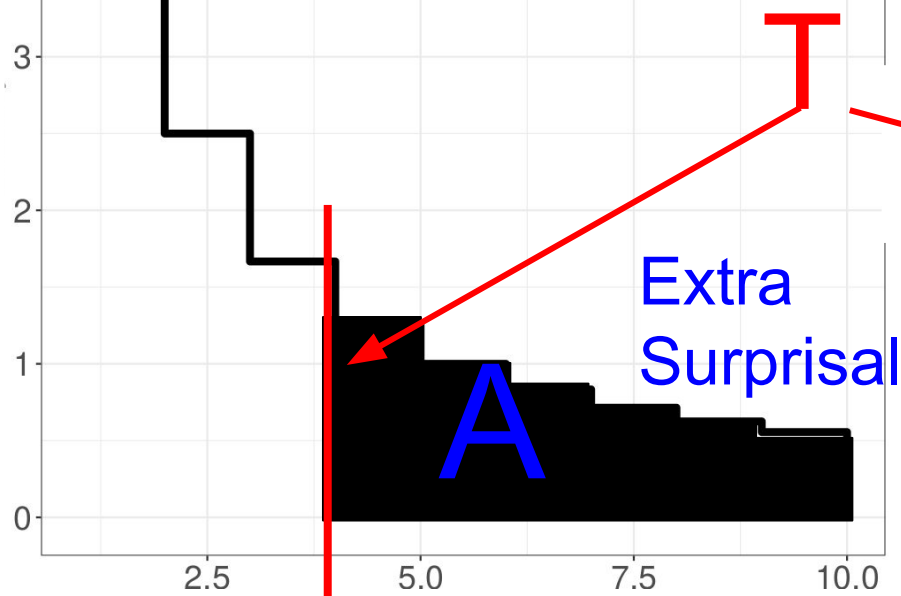


Intuition: Carrying information over long distances costs proportionally more.

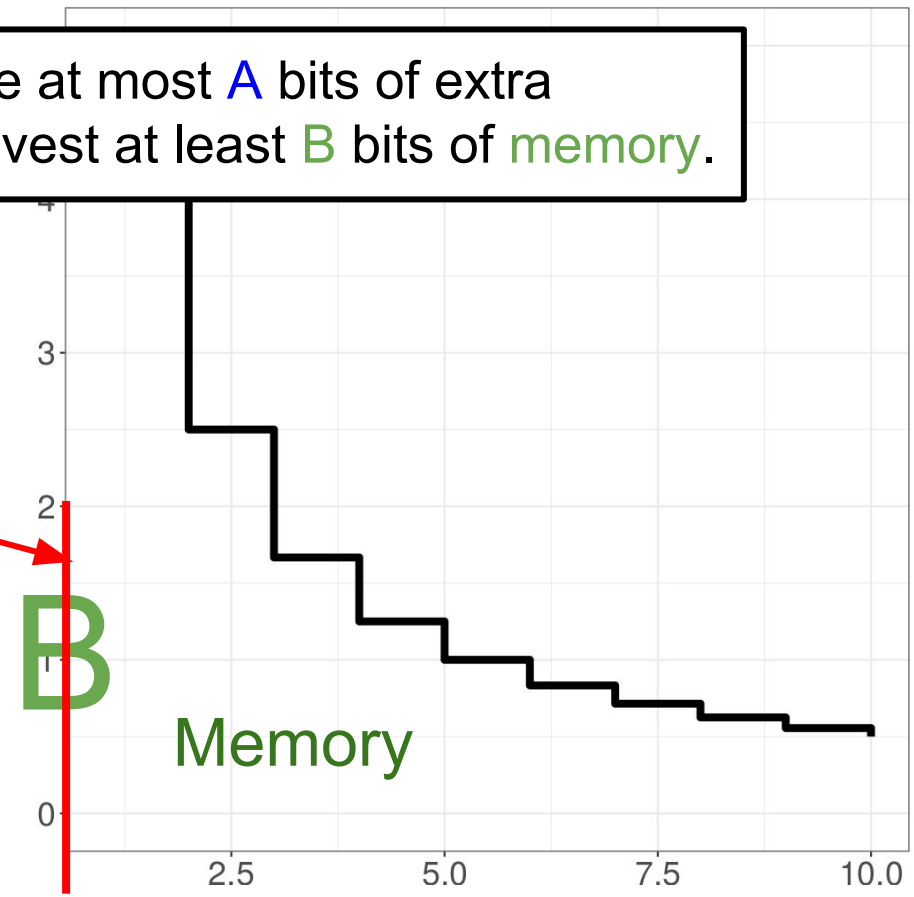
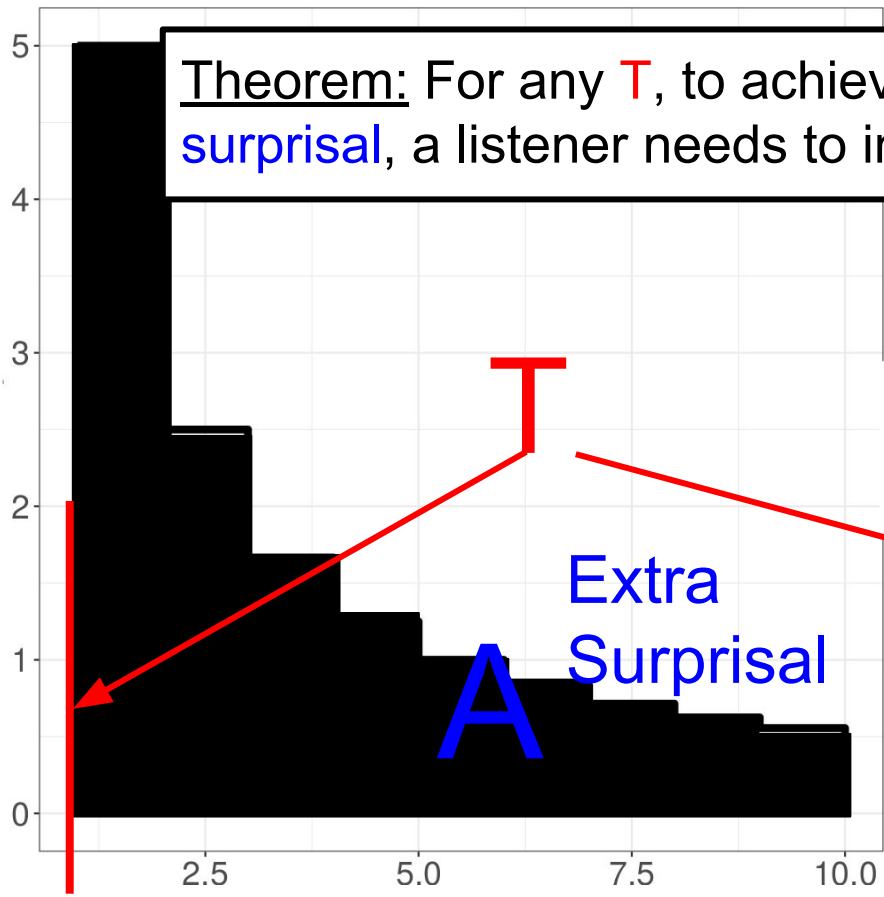


Theorem: For any T , to achieve at most A bits of extra surprisal, a listener needs to invest at least B bits of memory.

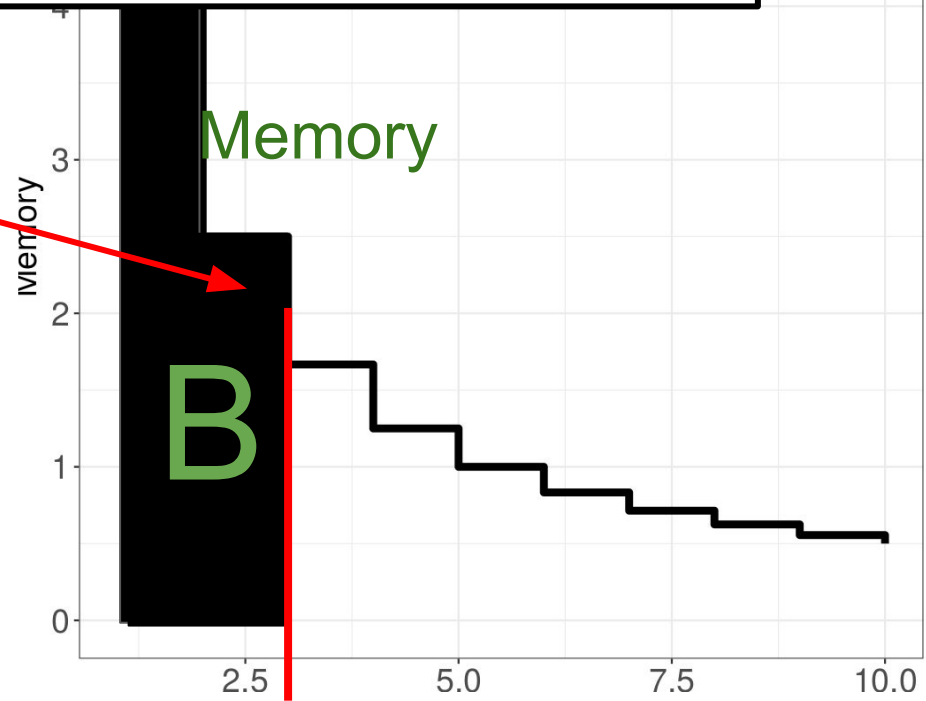
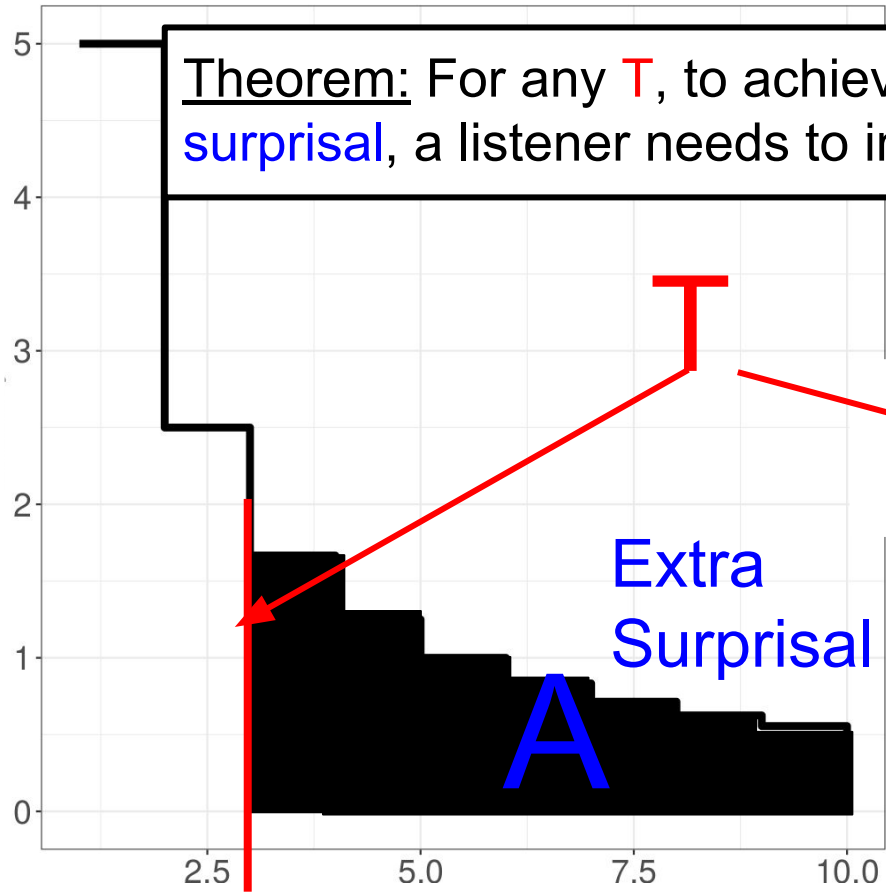
(Proof: bonus slides)



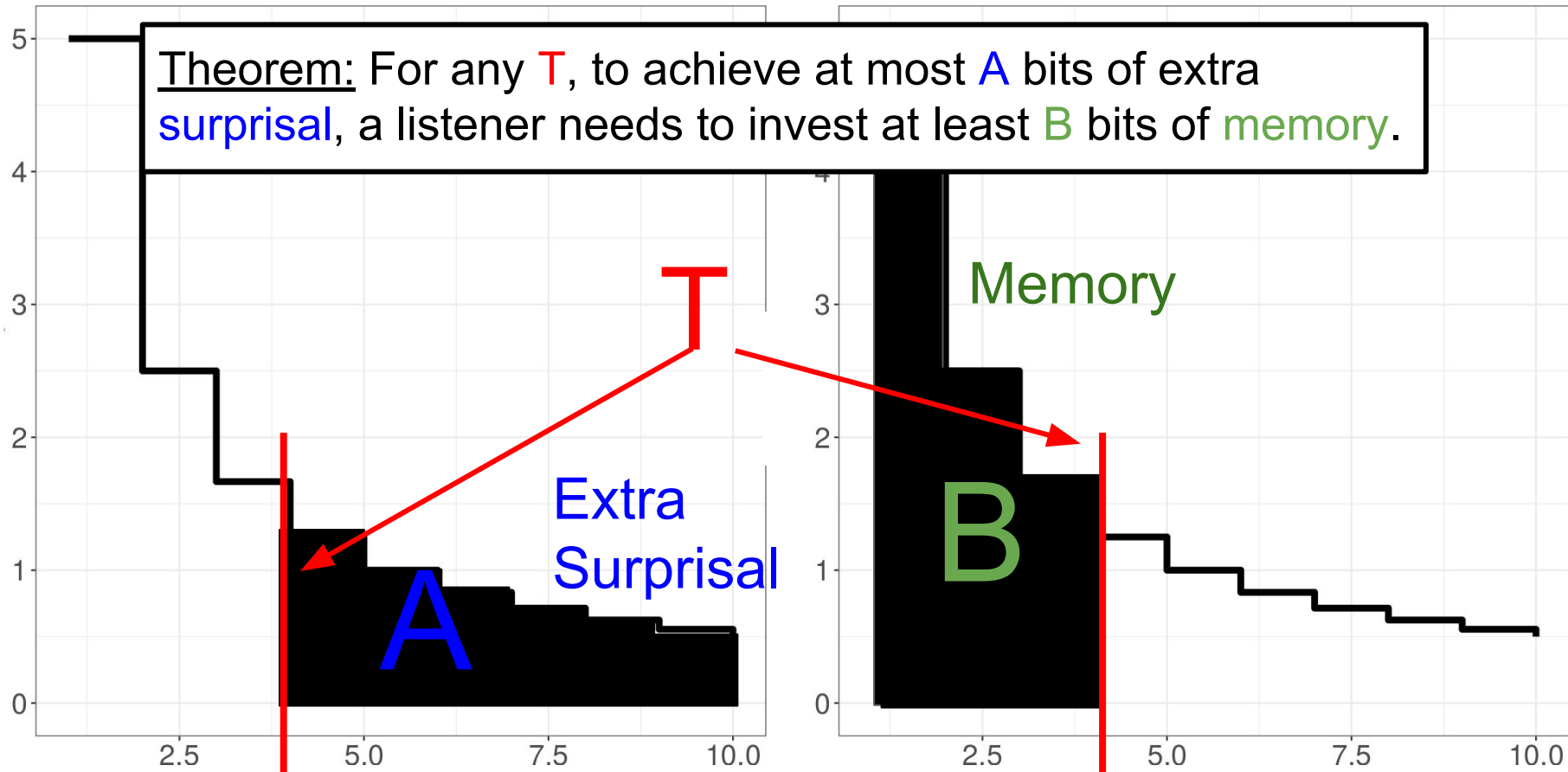
Theorem: For any T , to achieve at most A bits of extra surprisal, a listener needs to invest at least B bits of memory.



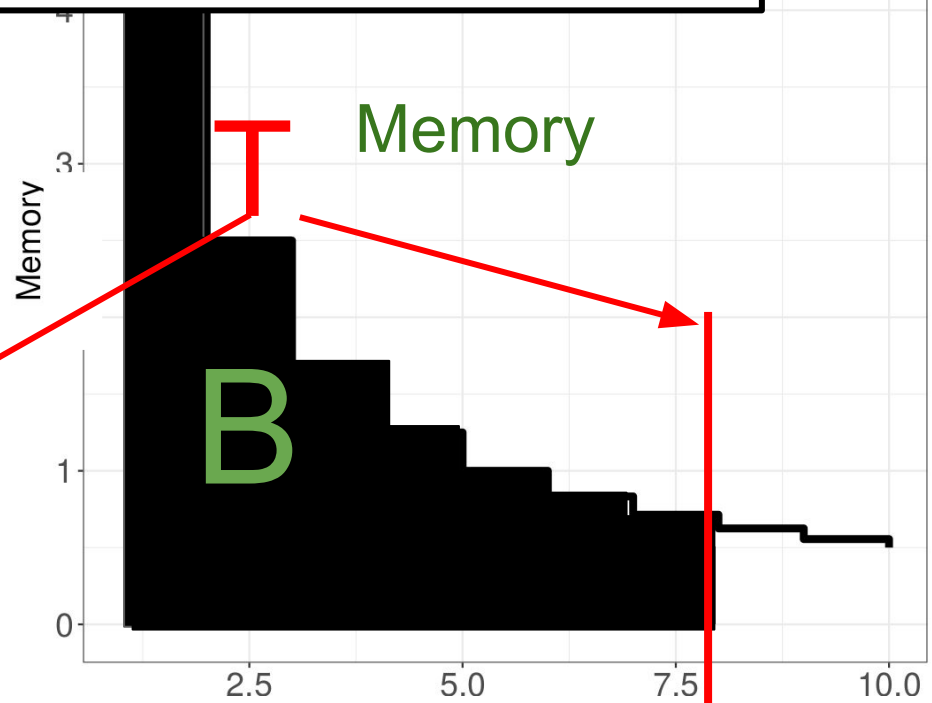
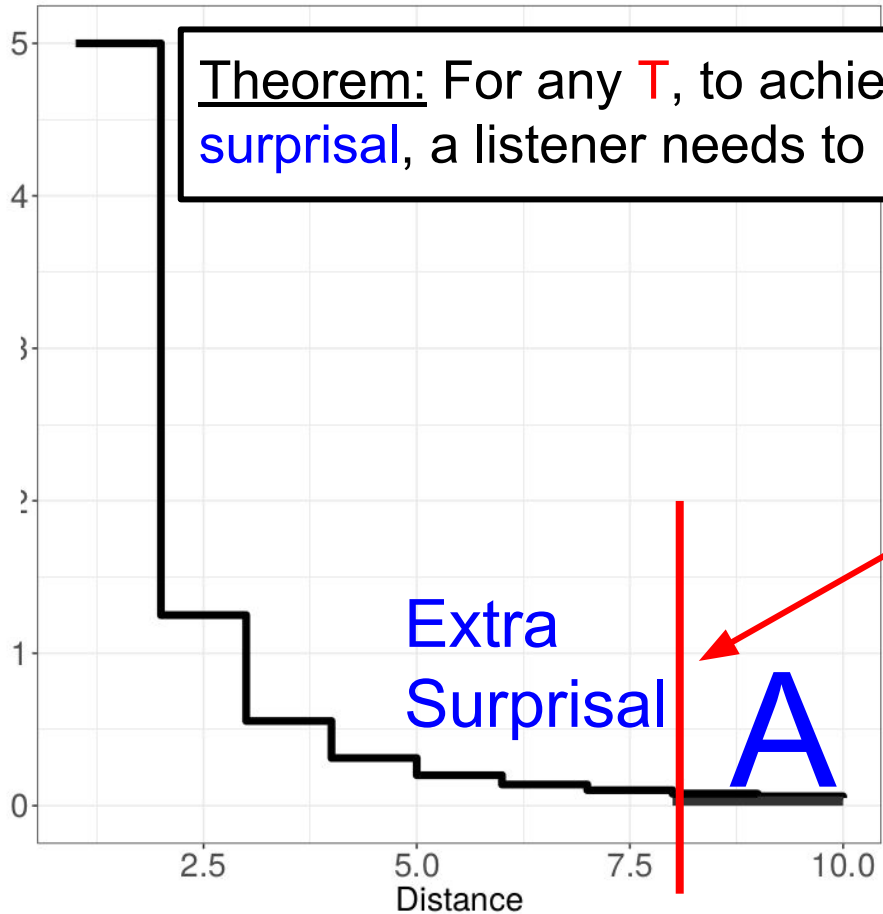
Theorem: For any T , to achieve at most A bits of extra surprisal, a listener needs to invest at least B bits of memory.



Theorem: For any T , to achieve at most A bits of extra surprisal, a listener needs to invest at least B bits of memory.



Theorem: For any T , to achieve at most A bits of extra surprisal, a listener needs to invest at least B bits of memory.

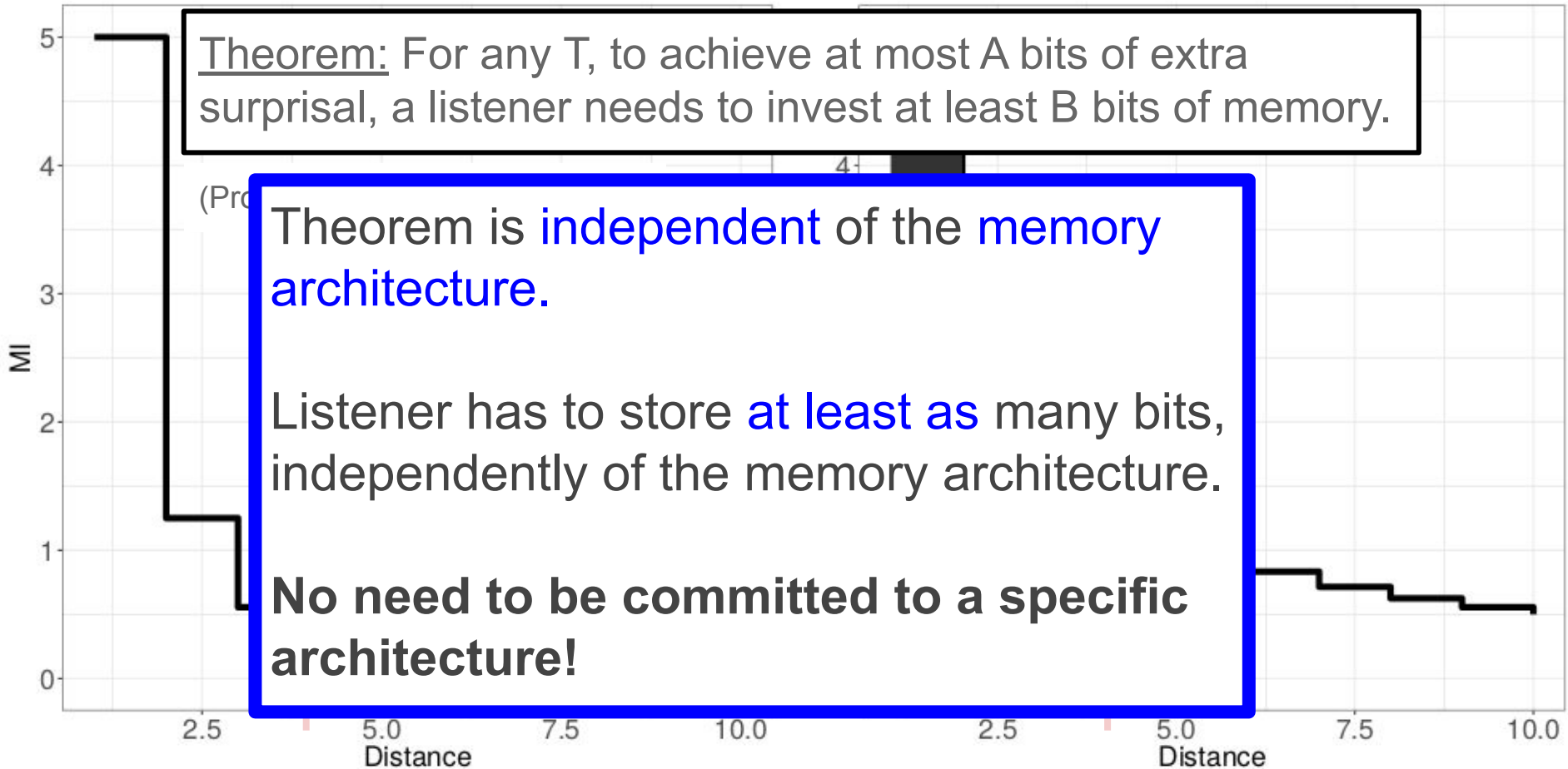


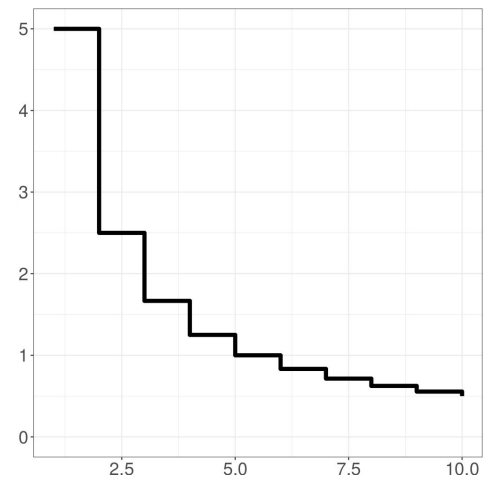
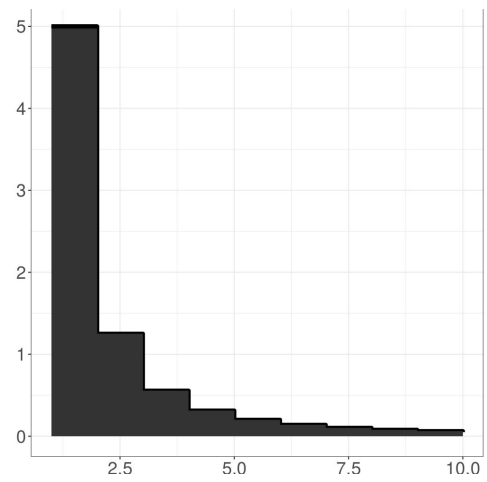
Theorem: For any T , to achieve at most A bits of extra surprisal, a listener needs to invest at least B bits of memory.

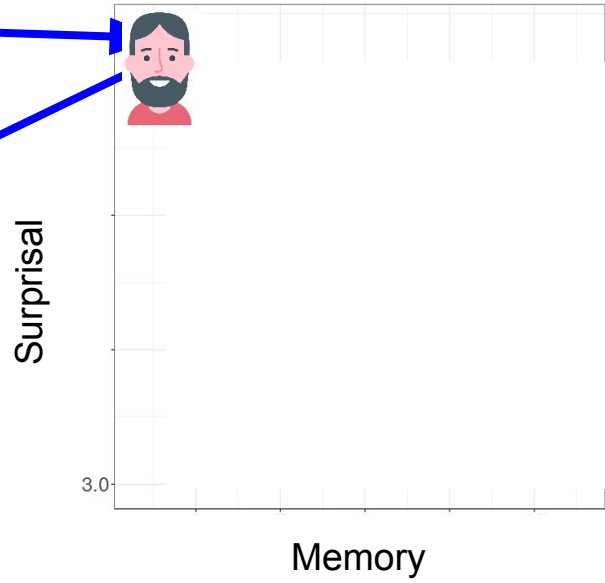
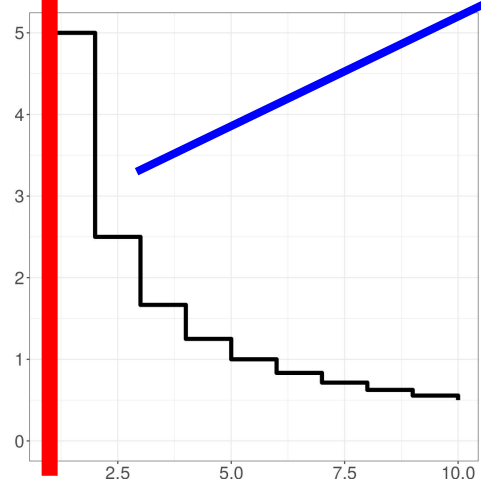
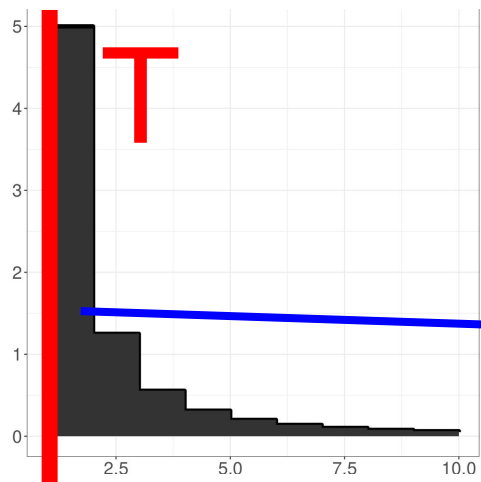
Theorem is **independent** of the **memory architecture**.

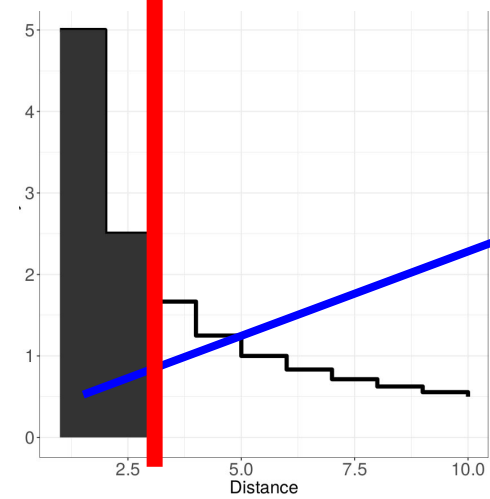
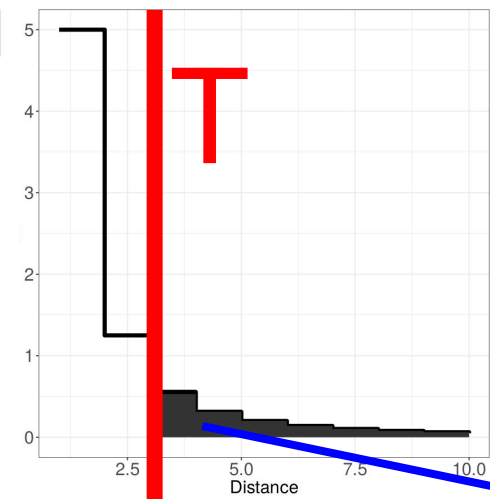
Listener has to store **at least as** many bits, independently of the memory architecture.

No need to be committed to a specific architecture!

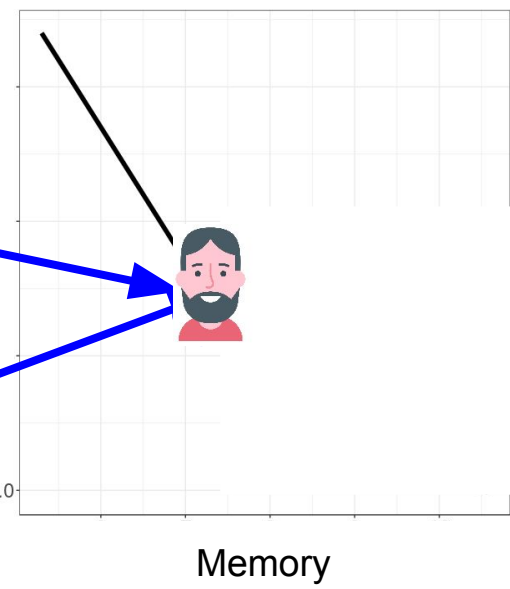


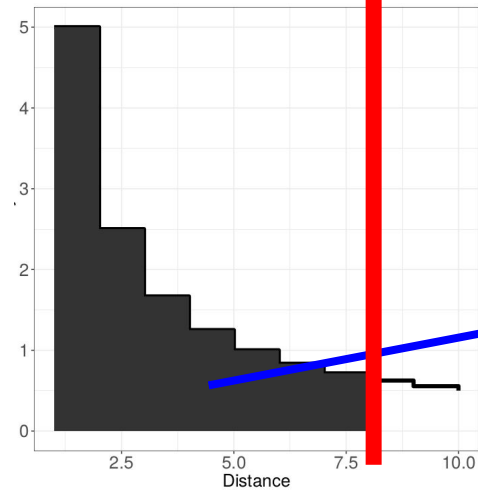
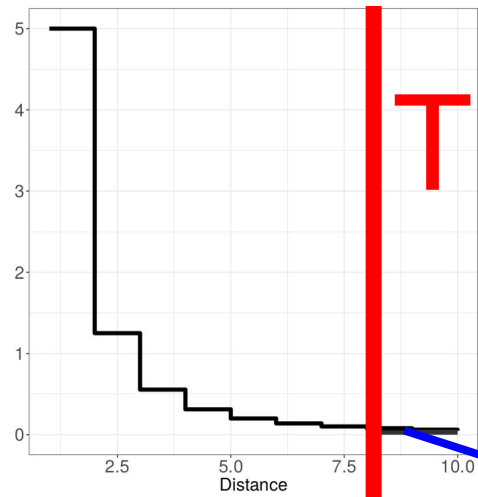




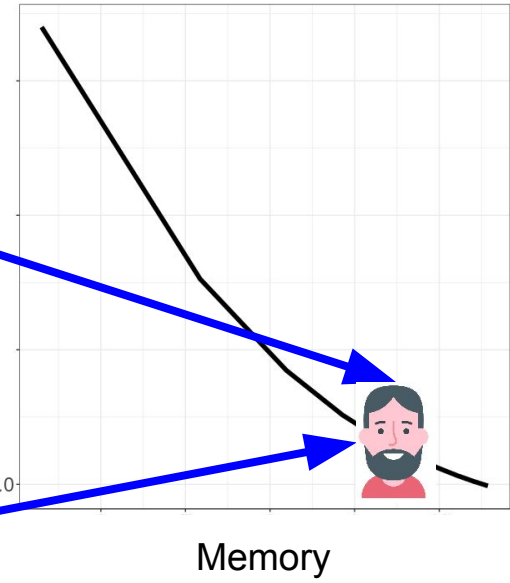


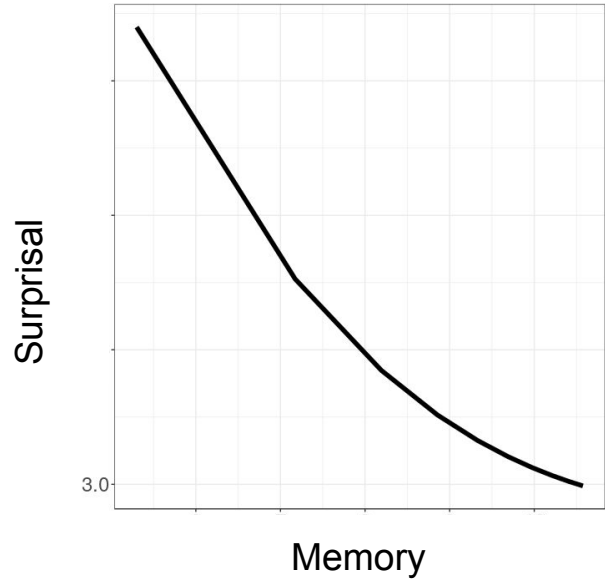
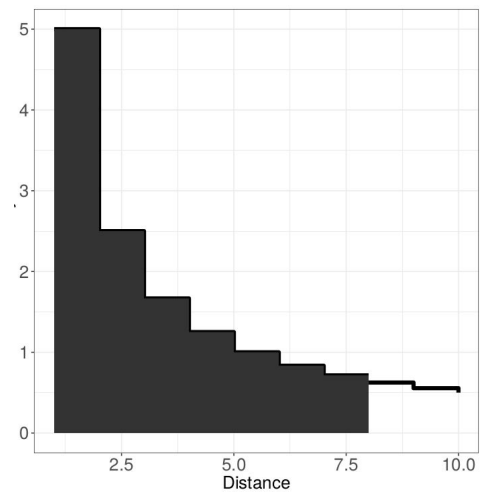
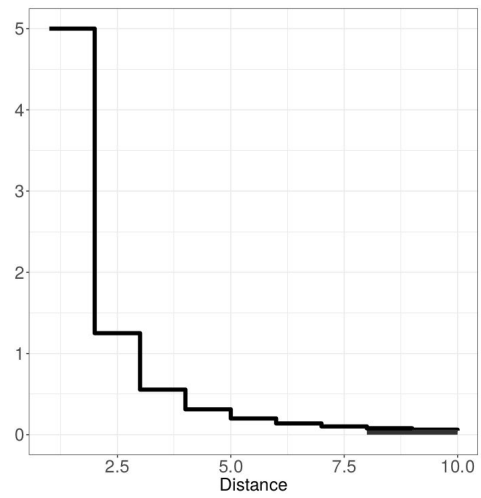
Surprisal





Surprisal





This talk

1. Information-theoretic formalization of memory limitations
2. Prove **theorem** describing **tradeoff between memory and surprisal**, without assumptions about memory architecture
3. Test: Are crosslinguistic word orders optimized for the memory-surprisal tradeoff?

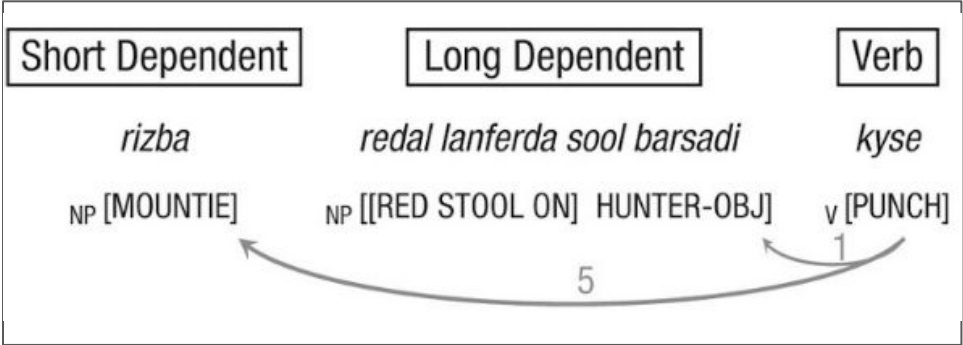
This talk

1. Information-theoretic formalization of memory limitations
2. Prove theorem describing tradeoff between memory and surprisal, without assumptions about memory architecture
3. Test: Are **crosslinguistic word orders optimized** for the memory-surprisal tradeoff?

Experiment 1:
Dependency Length in an Artificial
Language

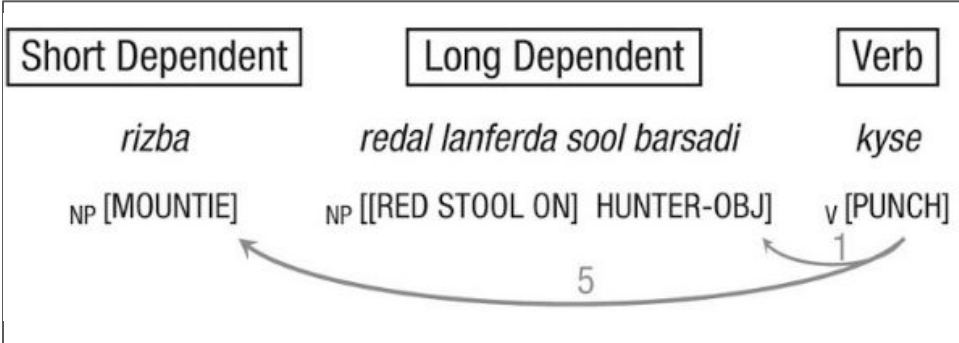
Dependency Length in an Artificial Language

Language A (long dependencies)

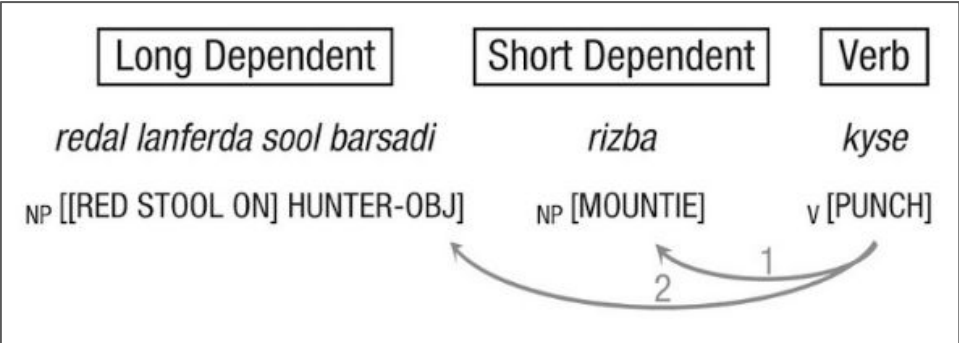


Dependency Length in an Artificial Language

Language A (long dependencies)

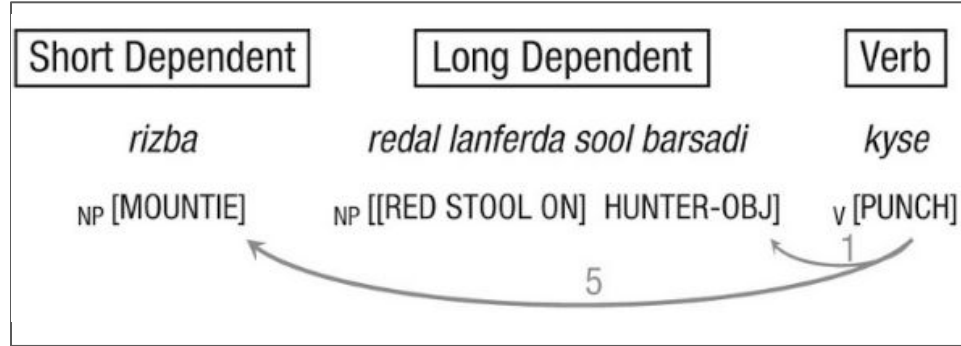


Language B (short dependencies)



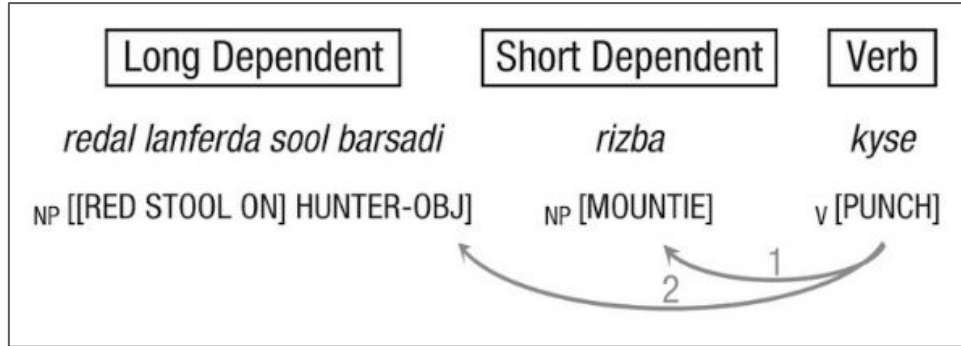
Dependency Length in an Artificial Language

Language A (long dependencies)



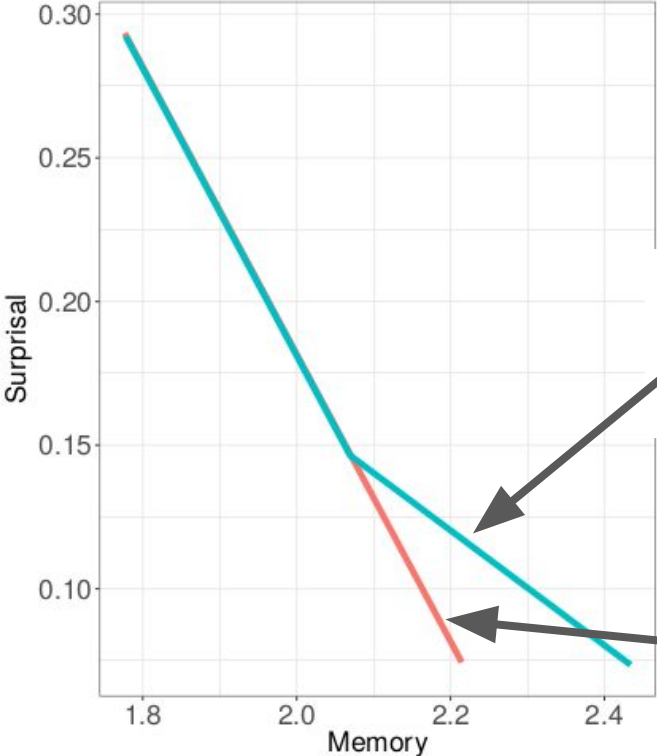
Participants tended to produce orders with **shorter dependencies**

Language B (short dependencies)

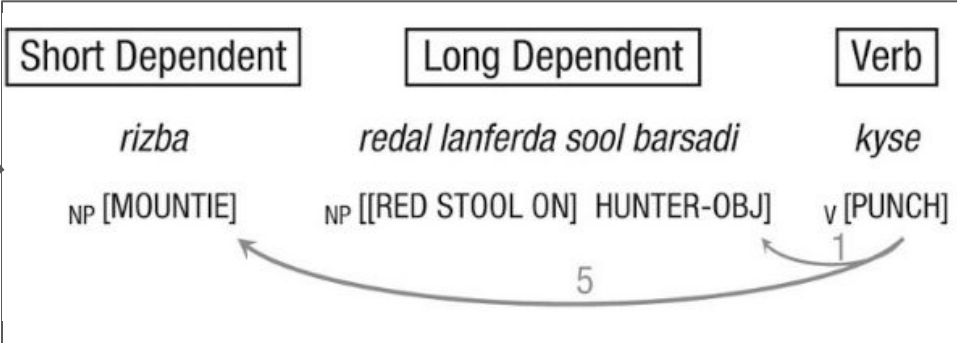


Fedzechkina et al. 2018

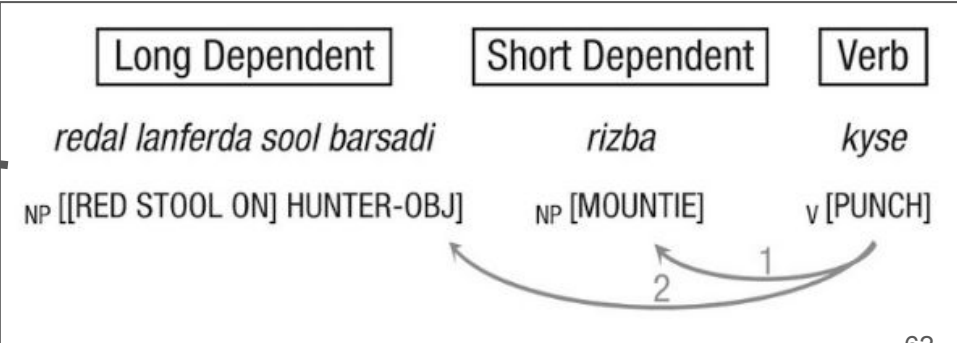
Dependency Length in an Artificial Language



Language A (long dependencies)



Language B (short dependencies)



Experiment 2: Crosslinguistic Word Orders

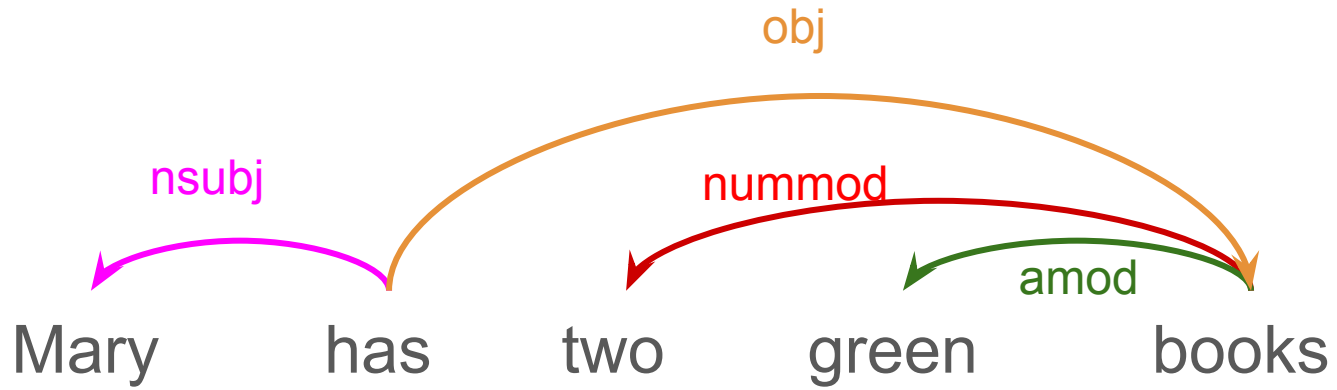
Question: Does language optimize the Memory-Surprisal tradeoff?

Method

1. **Syntactic corpora** from the Universal Dependencies Project (54 languages)
2. Create **counterfactual orderings** of the syntactic trees
3. Estimate **memory-surprisal tradeoff**
4. **Compare** memory need between **real** and **counterfactual** versions.

Method

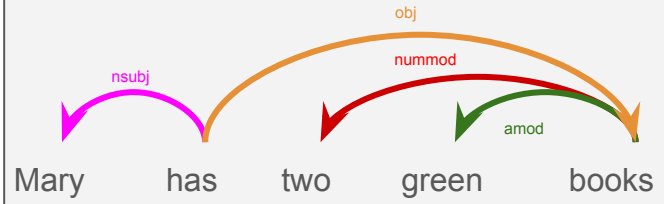
1. **Syntactic corpora from the Universal Dependencies Project (54 languages)**
2. Create counterfactual orderings of the syntactic trees
3. Estimate memory-surprisal tradeoff
4. Compare memory need between real and counterfactual versions.

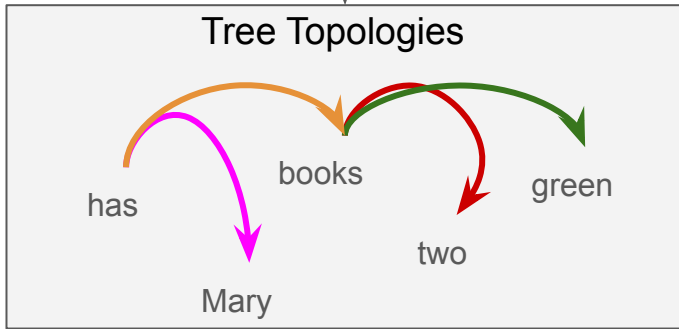
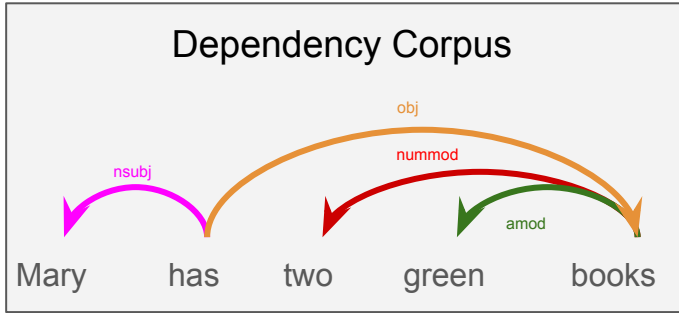


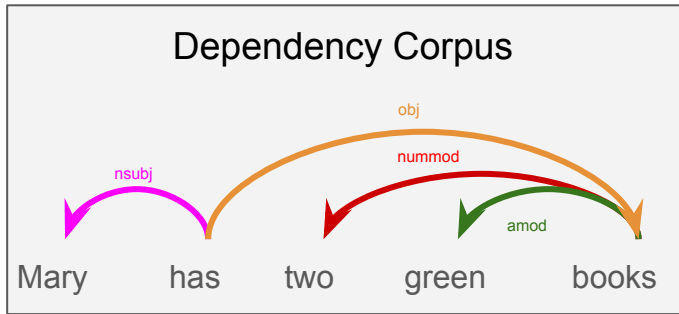
Method

1. Syntactic corpora from the Universal Dependencies Project (54 languages)
2. **Create counterfactual orderings of the syntactic trees**
3. Estimate memory-surprisal tradeoff
4. Compare memory need between real and counterfactual versions.

Dependency Corpus

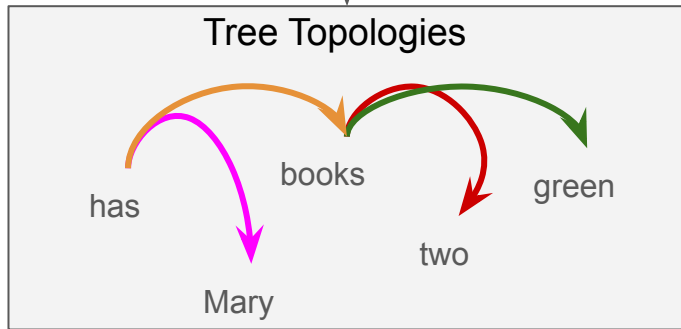


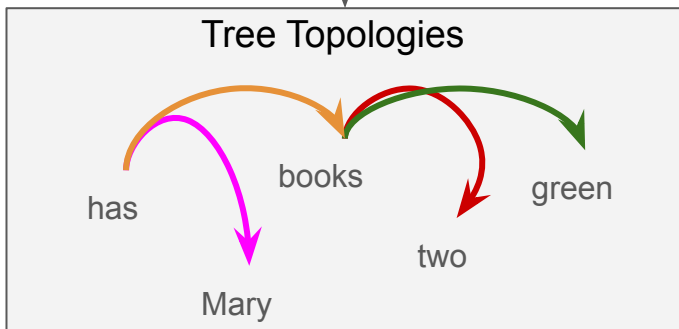
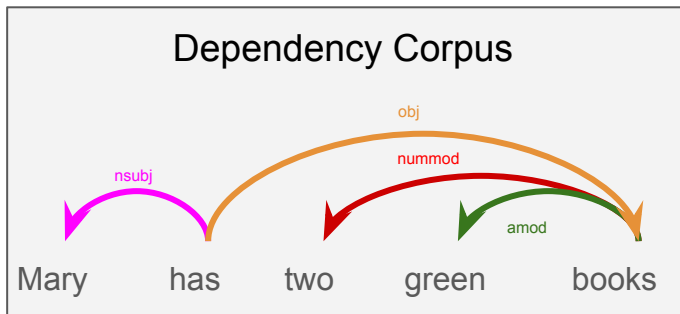




Ordering Grammar

NOUN	$\xrightarrow{\text{amod}}$	ADJ	0.3
NOUN	$\xrightarrow{\text{nummod}}$	NUM	0.7
VERB	$\xrightarrow{\text{nsubj}}$	NOUN	-0.2
VERB	$\xrightarrow{\text{obj}}$	NOUN	0.8
...			

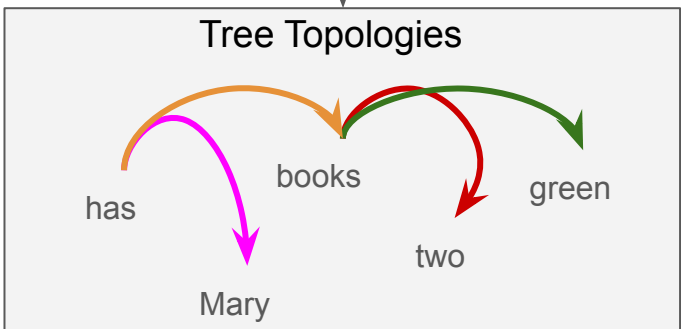
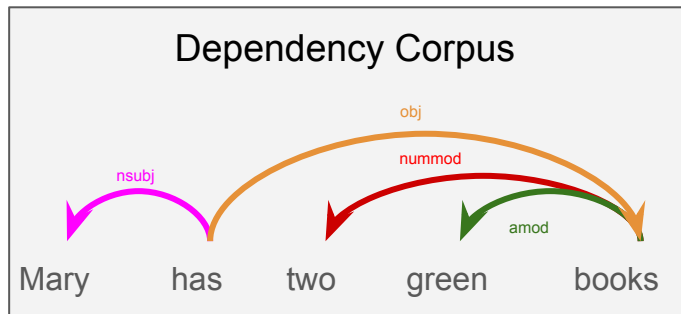




Ordering Grammar

NOUN	→ ^{amod}	ADJ	0.3
NOUN	→ ^{nummod}	NUM	0.7
VERB	→ ^{nsubj}	NOUN	-0.2
VERB	→ ^{obj}	NOUN	0.8
...			

“Object follows verb”

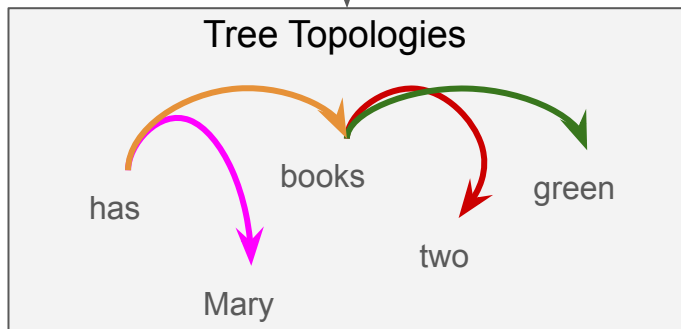
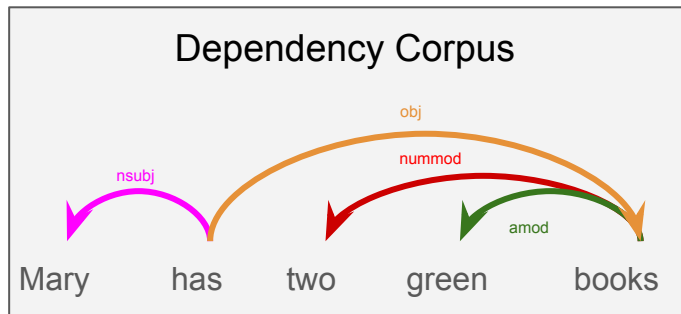


Ordering Grammar

NOUN	→ ^{amod}	ADJ	0.3
NOUN	→ ^{nummod}	NUM	0.7
VERB	→ ^{nsubj}	NOUN	-0.2
VERB	→ ^{obj}	NOUN	0.8
...			

“Object follows verb”

“Adjective precedes noun”



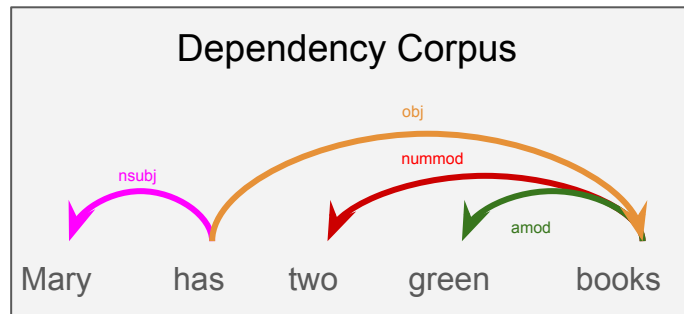
Ordering Grammar

NOUN	→ ^{amod}	ADJ	0.3
NOUN	→ ^{nummod}	NUM	0.7
VERB	→ ^{nsubj}	NOUN	-0.2
VERB	→ ^{obj}	NOUN	0.8
...			

“Object follows verb”

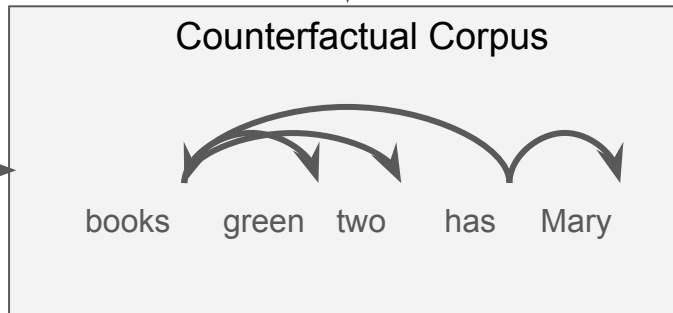
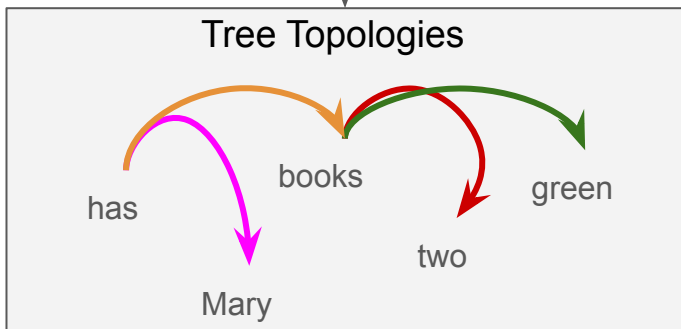
“Adjective precedes noun”

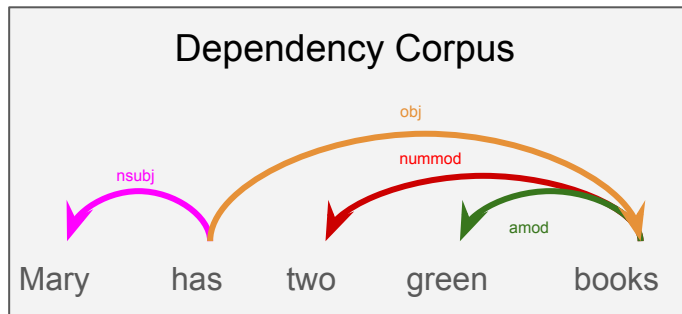
“Numerals follow adjectives & precede nouns”



Ordering Grammar

NOUN	→	ADJ	0.3
NOUN	→	NUM	0.7
VERB	→	NOUN	-0.2
VERB	→	NOUN	0.8
...			

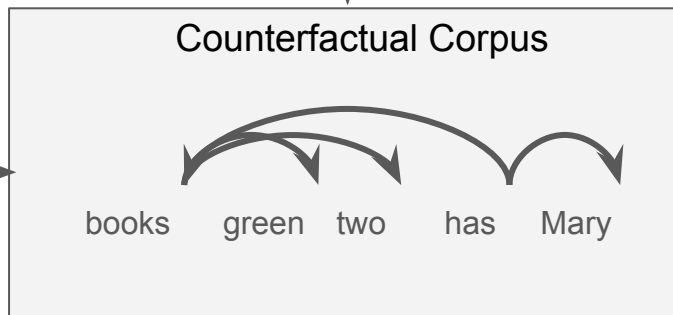
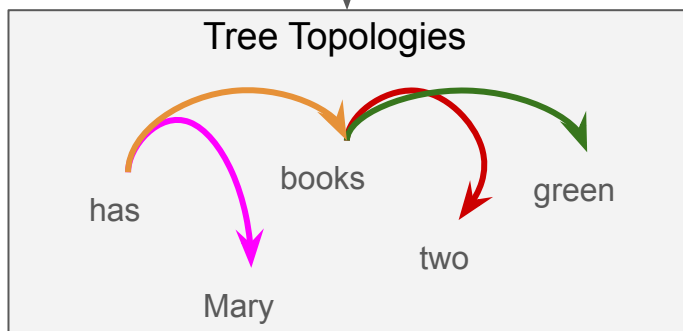


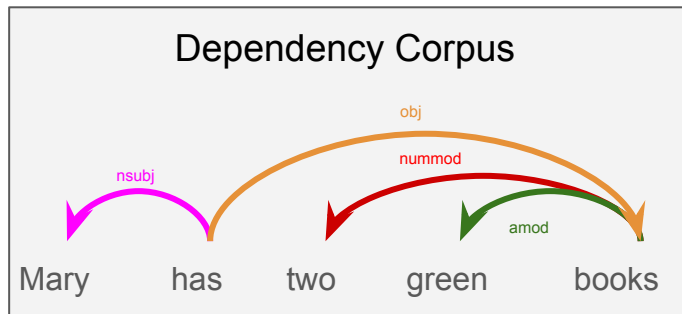


Ordering Grammar

NOUN	$\xrightarrow{\text{amod}}$	ADJ	0.3
NOUN	$\xrightarrow{\text{nummod}}$	NUM	0.7
VERB	$\xrightarrow{\text{nsubj}}$	NOUN	-0.2
VERB	$\xrightarrow{\text{obj}}$	NOUN	0.8
...			

Each **parameter setting** generates a **different counterfactual corpus.**

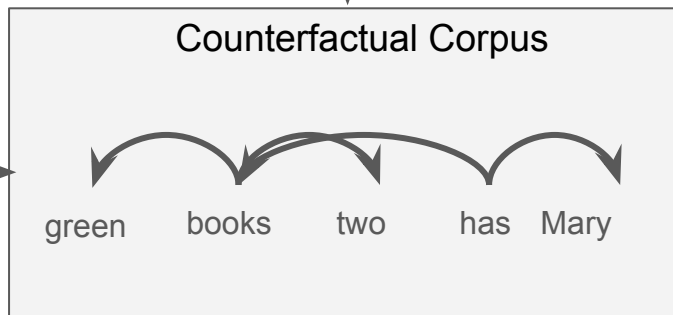
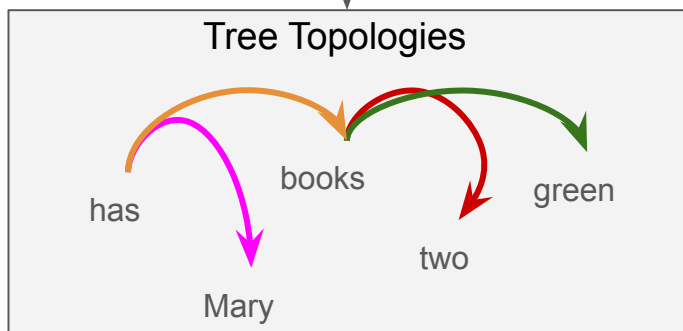


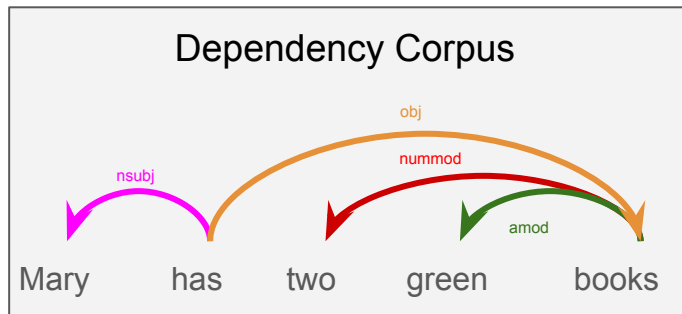


Ordering Grammar

NOUN	$\xrightarrow{\text{amod}}$	ADJ	0.9
NOUN	$\xrightarrow{\text{nummod}}$	NUM	0.1
VERB	$\xrightarrow{\text{nsubj}}$	NOUN	0.5
VERB	$\xrightarrow{\text{obj}}$	NOUN	0.2
...			

Each **parameter setting** generates a **different counterfactual corpus.**

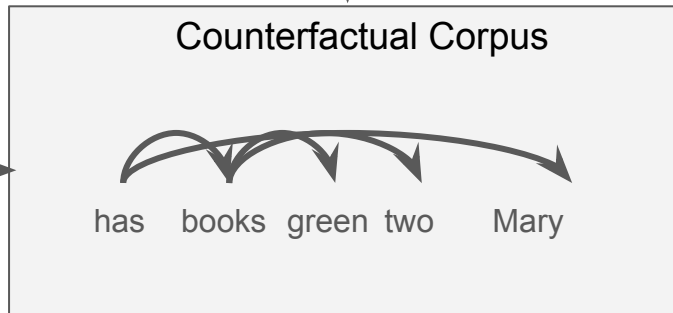
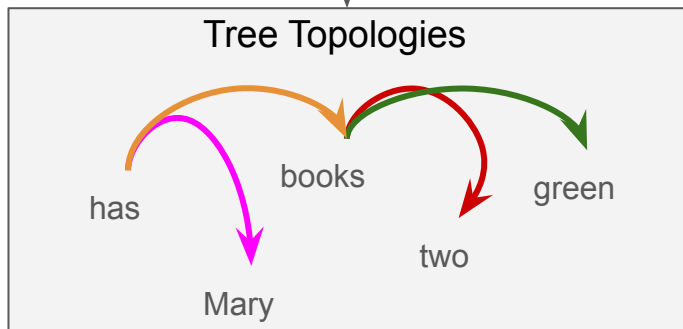


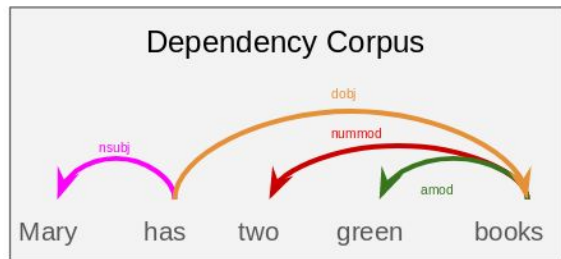


Ordering Grammar

NOUN	$\xrightarrow{\text{amod}}$	ADJ	0.1
NOUN	$\xrightarrow{\text{nummod}}$	NUM	0.95
VERB	$\xrightarrow{\text{nsubj}}$	NOUN	0.42
VERB	$\xrightarrow{\text{obj}}$	NOUN	0.82
...			

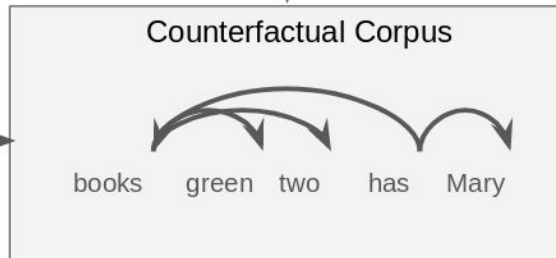
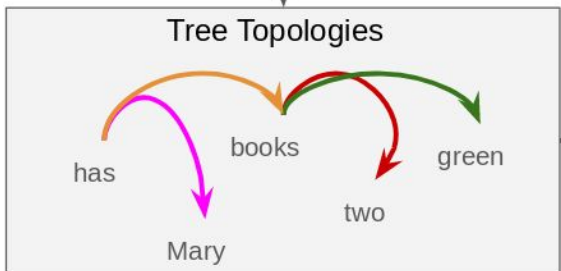
Each **parameter setting** generates a **different counterfactual corpus.**



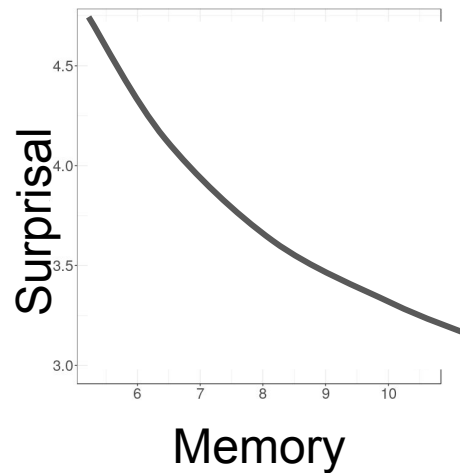


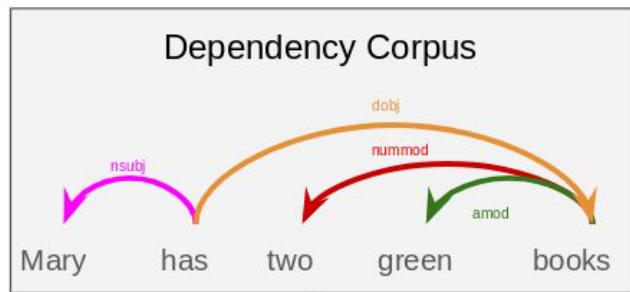
Ordering Grammar

NOUN	$\xrightarrow{\text{amod}}$	ADJ	0.3
NOUN	$\xrightarrow{\text{nummod}}$	NUM	0.7
VERB	$\xrightarrow{\text{nsubj}}$	NOUN	-0.2
VERB	$\xrightarrow{\text{dobj}}$	NOUN	0.8
...			



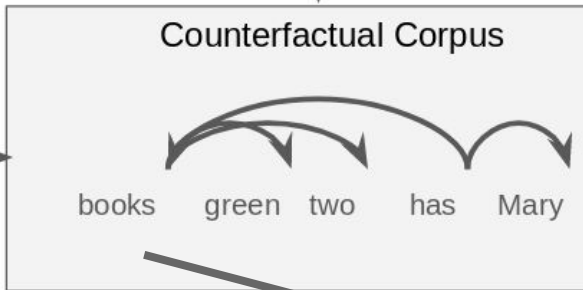
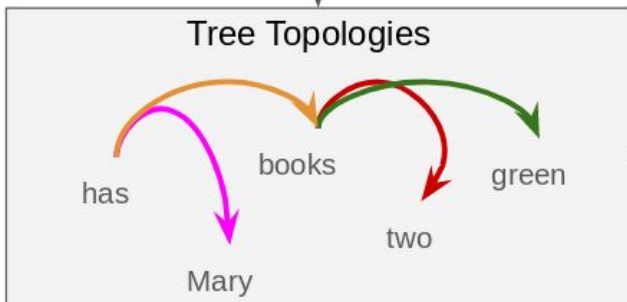
We compute **memory-surprisal tradeoff** on counterfactual corpora.



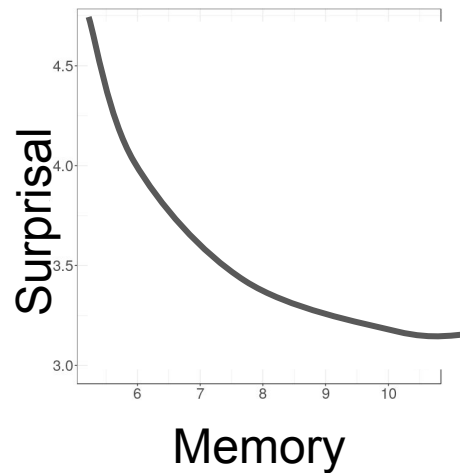


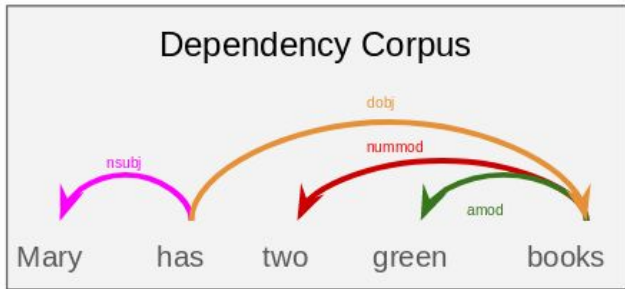
Ordering Grammar

NOUN	$\xrightarrow{\text{amod}}$	ADJ	0.3
NOUN	$\xrightarrow{\text{nummod}}$	NUM	0.7
VERB	$\xrightarrow{\text{nsubj}}$	NOUN	-0.2
VERB	$\xrightarrow{\text{dobj}}$	NOUN	0.8
...			



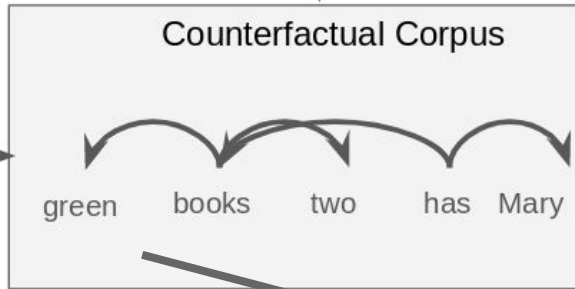
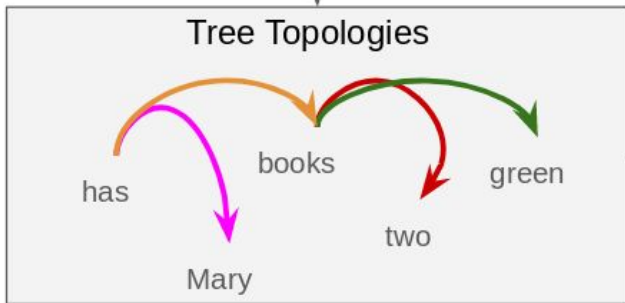
We compute **memory-surprisal tradeoff** on counterfactual corpora.



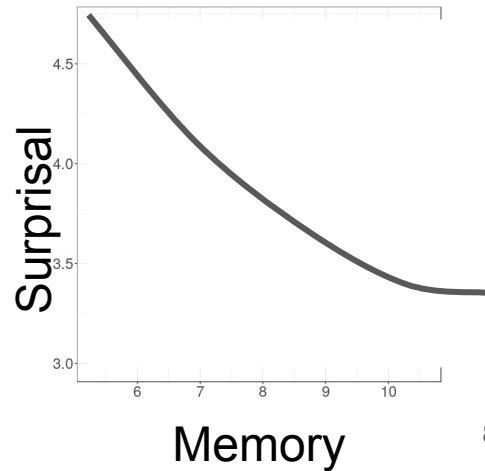


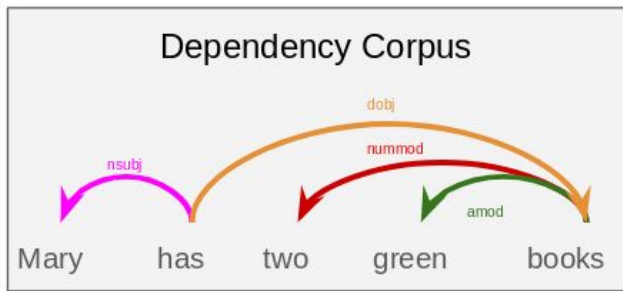
Ordering Grammar

NOUN	$\xrightarrow{\text{amod}}$	ADJ	0.9
NOUN	$\xrightarrow{\text{nummod}}$	NUM	0.1
VERB	$\xrightarrow{\text{nsubj}}$	NOUN	0.5
VERB	$\xrightarrow{\text{dobj}}$	NOUN	0.2
...			



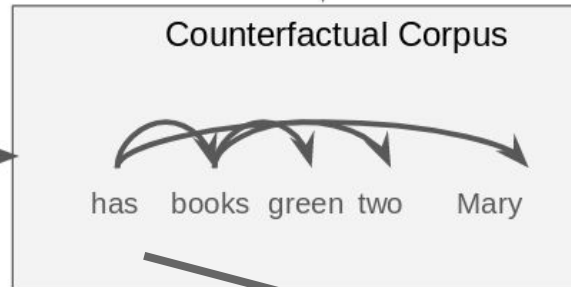
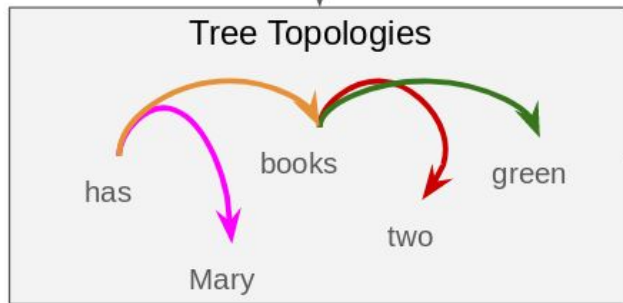
We compute **memory-surprisal tradeoff** on counterfactual corpora.



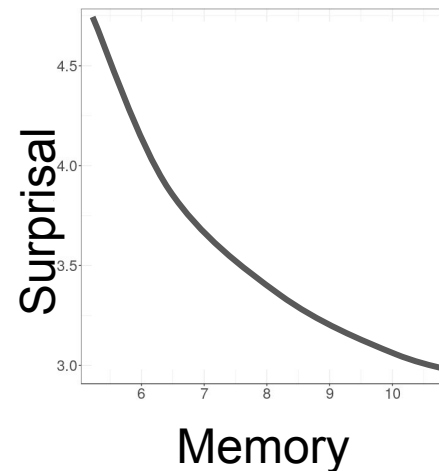


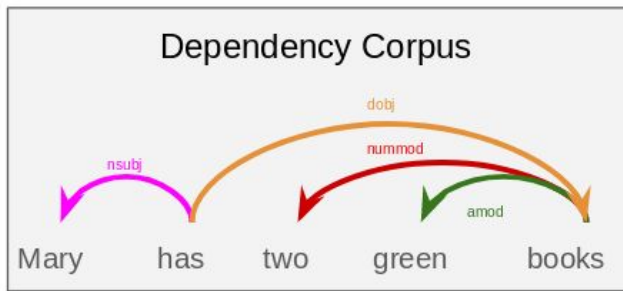
Ordering Grammar

NOUN	$\xrightarrow{\text{amod}}$	ADJ	0.1
NOUN	$\xrightarrow{\text{nummod}}$	NUM	0.95
VERB	$\xrightarrow{\text{nsubj}}$	NOUN	04.2
VERB	$\xrightarrow{\text{dobj}}$	NOUN	0.82
...			



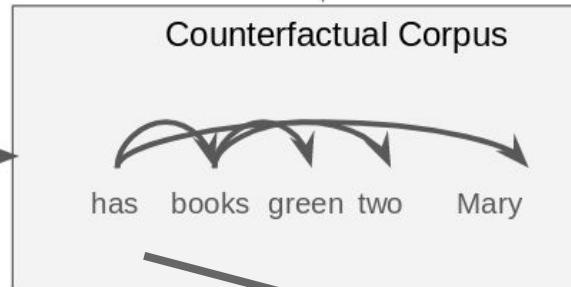
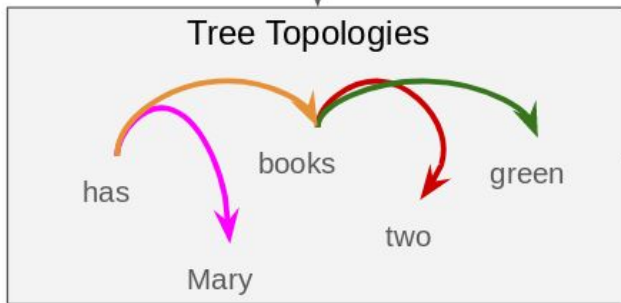
We compute **memory-surprisal tradeoff** on counterfactual corpora.



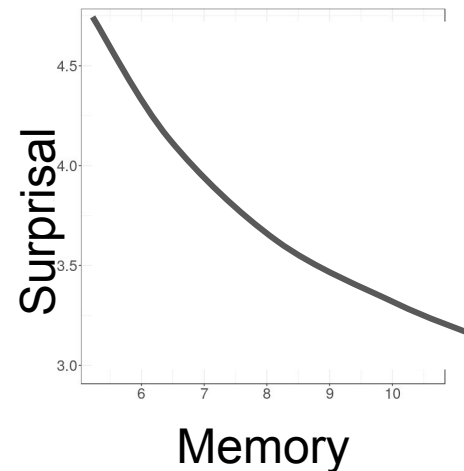


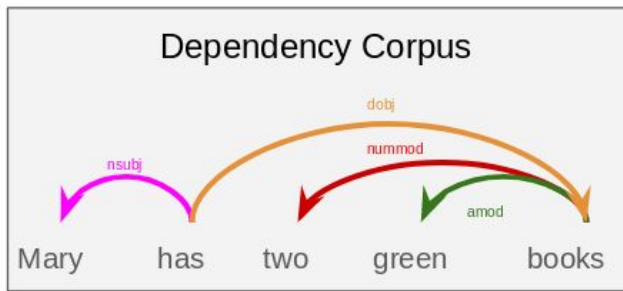
Ordering Grammar

NOUN	$\xrightarrow{\text{amod}}$	ADJ	0.1
NOUN	$\xrightarrow{\text{nummod}}$	NUM	0.95
VERB	$\xrightarrow{\text{nsubj}}$	NOUN	04.2
VERB	$\xrightarrow{\text{dobj}}$	NOUN	0.82
...			



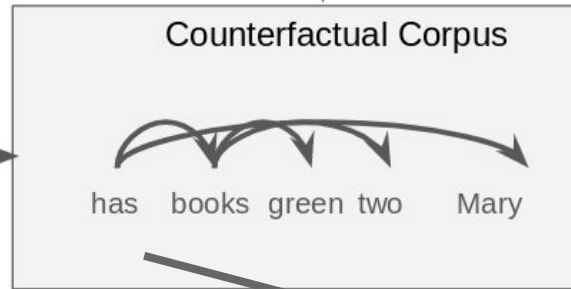
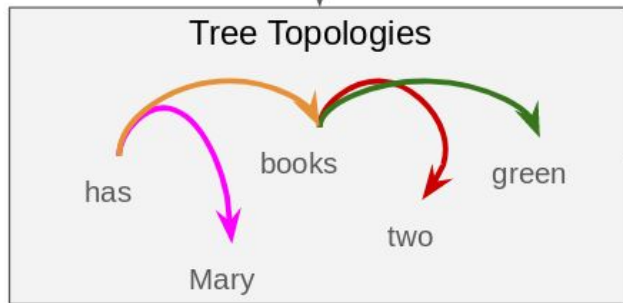
We compute **memory-surprisal tradeoff** on counterfactual corpora.



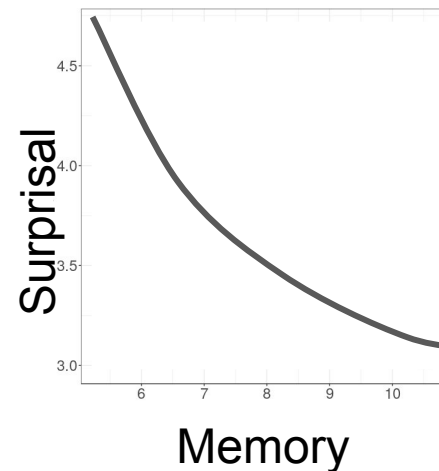


Ordering Grammar

NOUN	$\xrightarrow{\text{amod}}$	ADJ	0.1
NOUN	$\xrightarrow{\text{nummod}}$	NUM	0.95
VERB	$\xrightarrow{\text{nsubj}}$	NOUN	04.2
VERB	$\xrightarrow{\text{dobj}}$	NOUN	0.82
...			

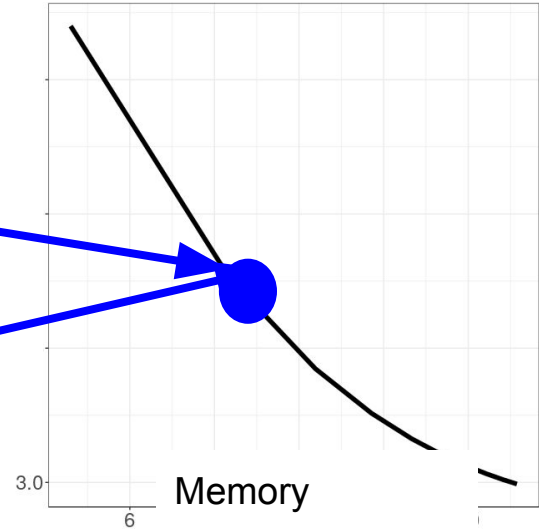
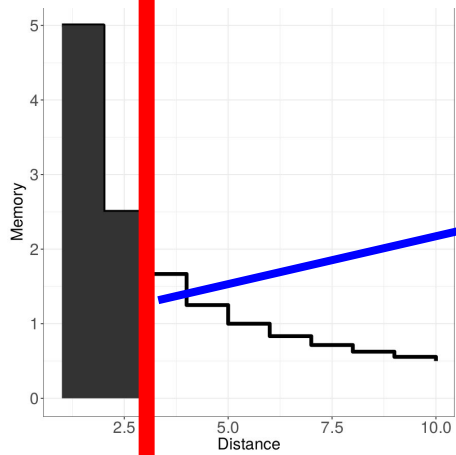
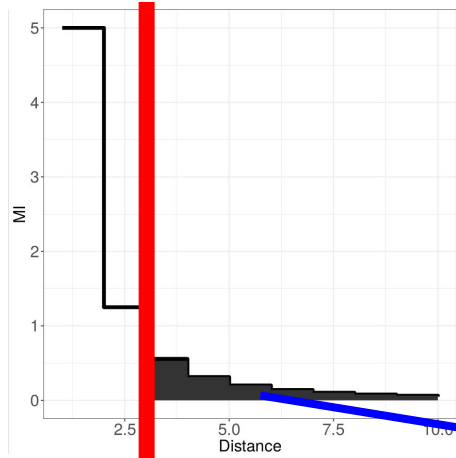


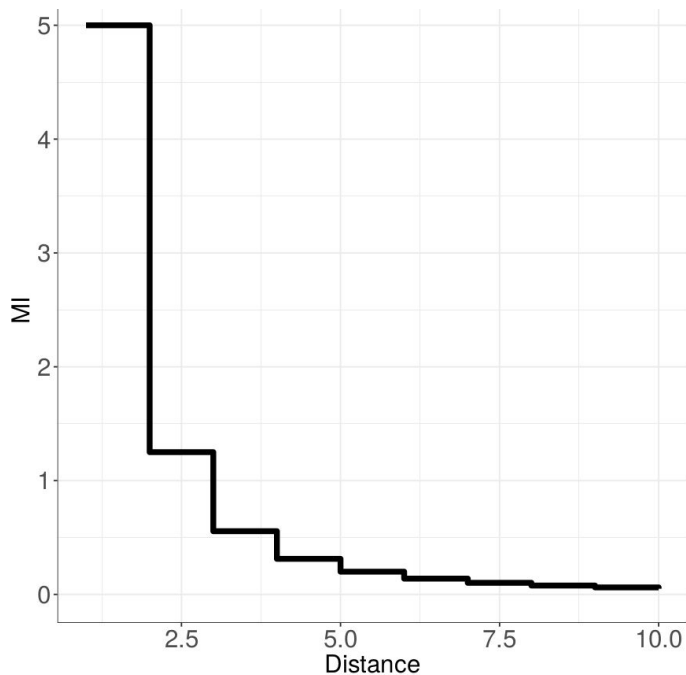
We compute **memory-surprisal tradeoff** on counterfactual corpora.



Method

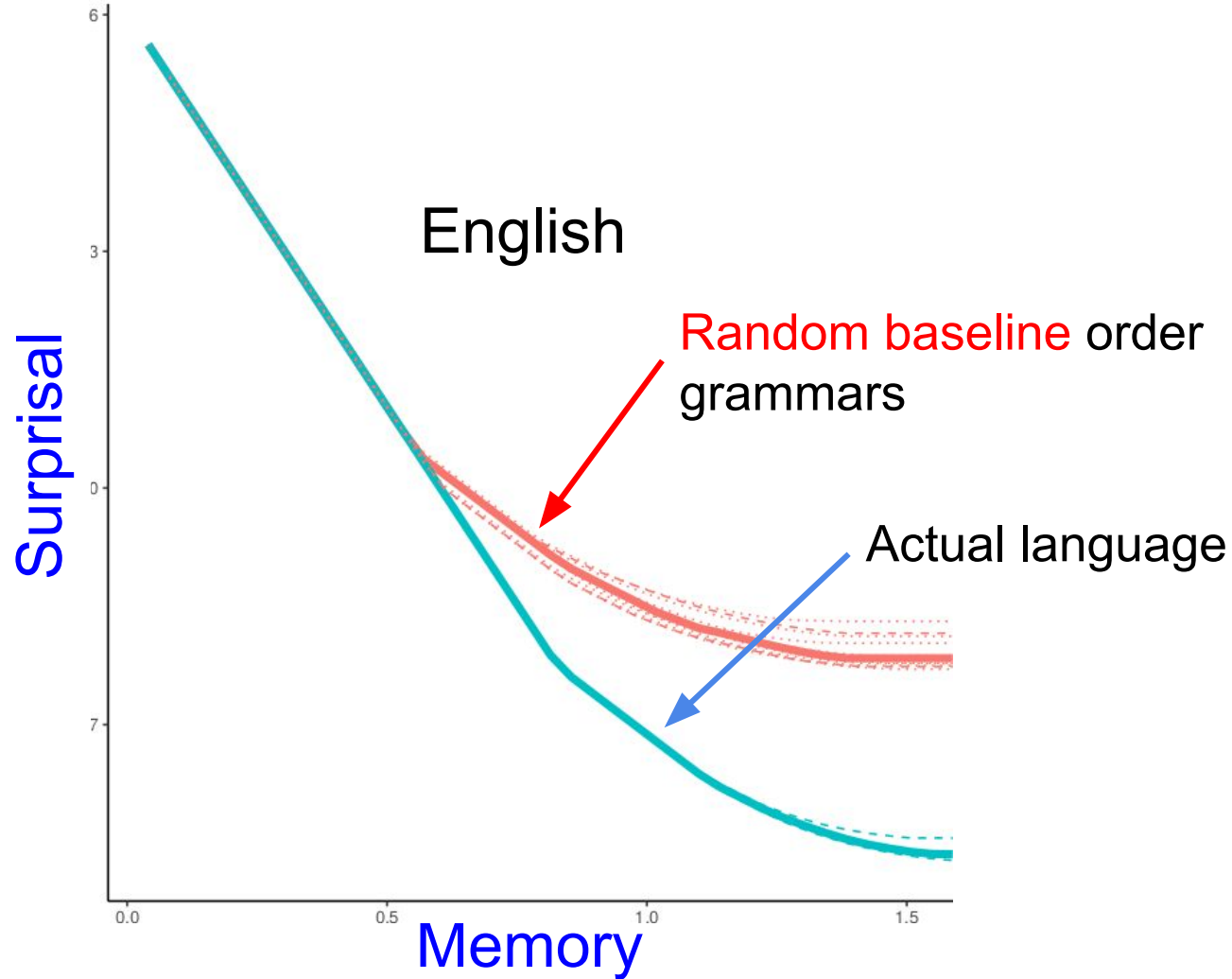
1. Syntactic corpora from the Universal Dependencies Project (54 languages)
2. Create counterfactual orderings of the syntactic trees
3. **Estimate memory-surprisal tradeoff**
4. Compare memory need between real and counterfactual versions.

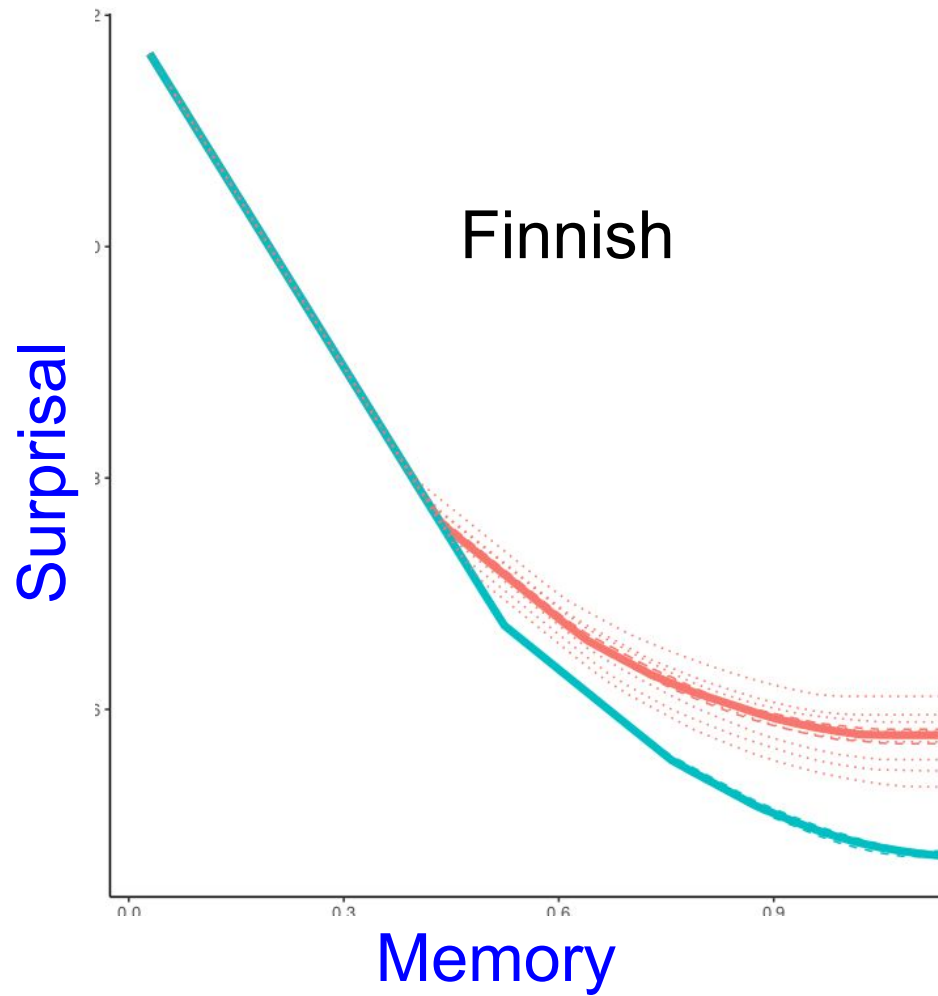


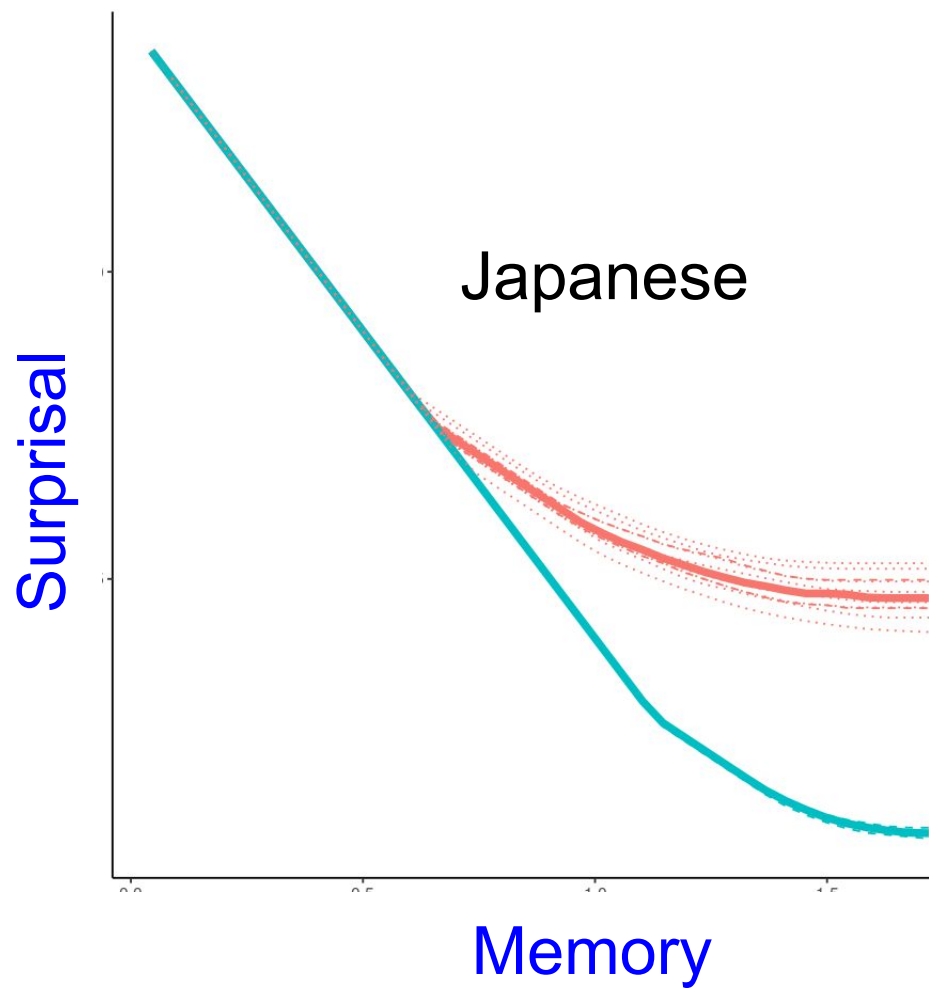


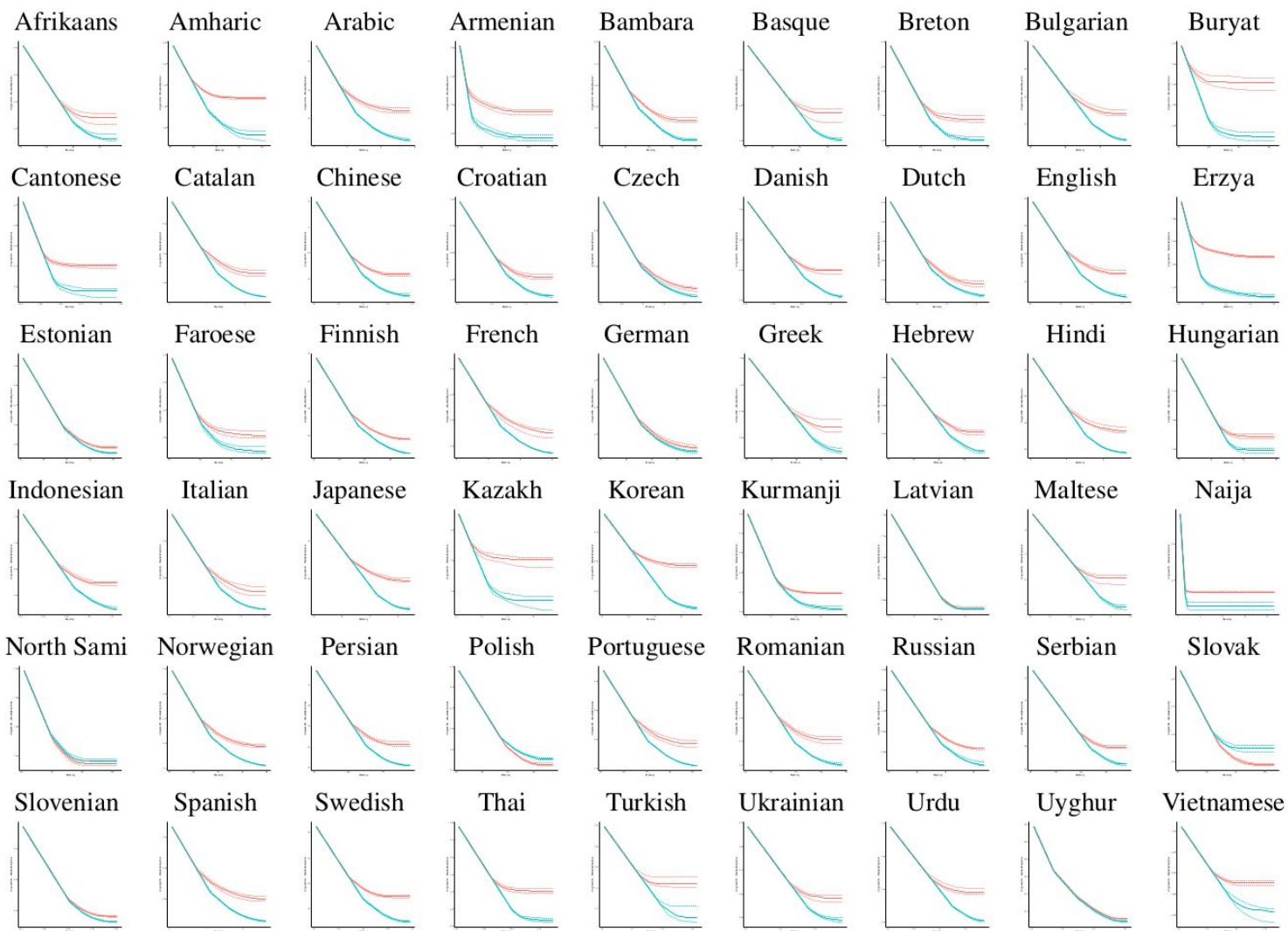
Estimated using LSTM recurrent neural networks

- essentially the state of the art in statistical modeling of language
- similar results obtained using traditional methods (transition probabilities & n-gram models)











Real orderings leads to better tradeoff ($p < 0.001$) in 50 out of 54 languages

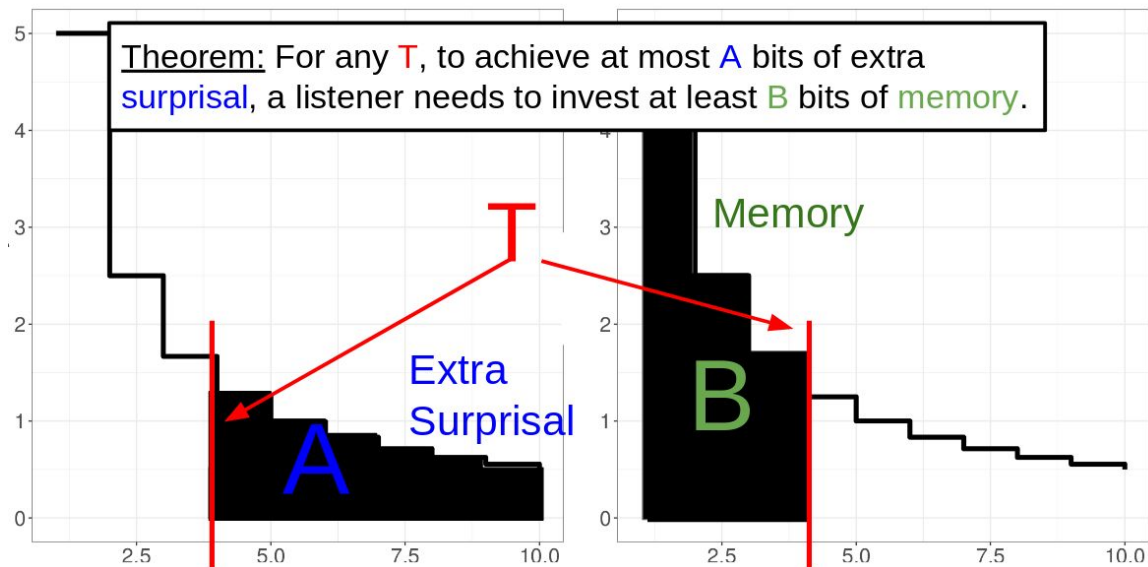
Conclusions

There is a tradeoff between **listener memory** and **experienced surprisal**.

Conclusions

There is a tradeoff between **listener memory** and **experienced surprisal**.

We formalize it using **Information Theory**, minimizing architectural assumptions

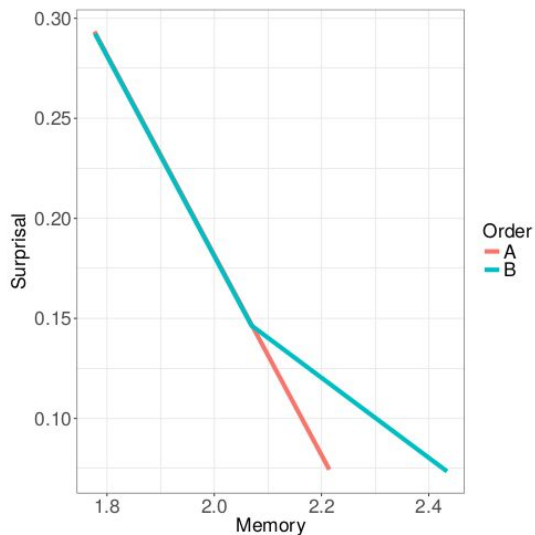


Conclusions

There is a tradeoff between **listener memory** and **experienced surprisal**.

We formalize it using **Information Theory**, minimizing architectural assumptions

Languages with **short dependencies** have better tradeoffs.



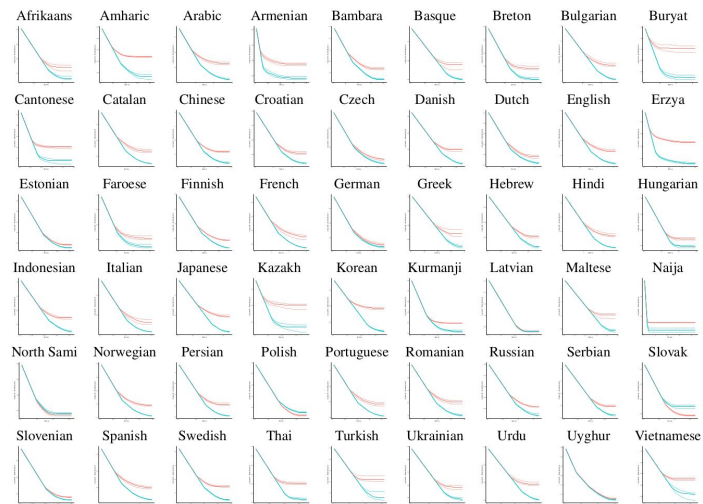
Conclusions

There is a tradeoff between **listener memory** and **experienced surprisal**.

We formalize it using **Information Theory**, minimizing architectural assumptions

Languages with **short dependencies** have better tradeoffs.

Crosslinguistic word orders support **more efficient tradeoffs** than most counterfactual orders.



Thanks!

Proof



X_{-6}

X_{-5}

X_{-4}

X_{-3}

X_{-2}

X_{-1}

X_1

X_2

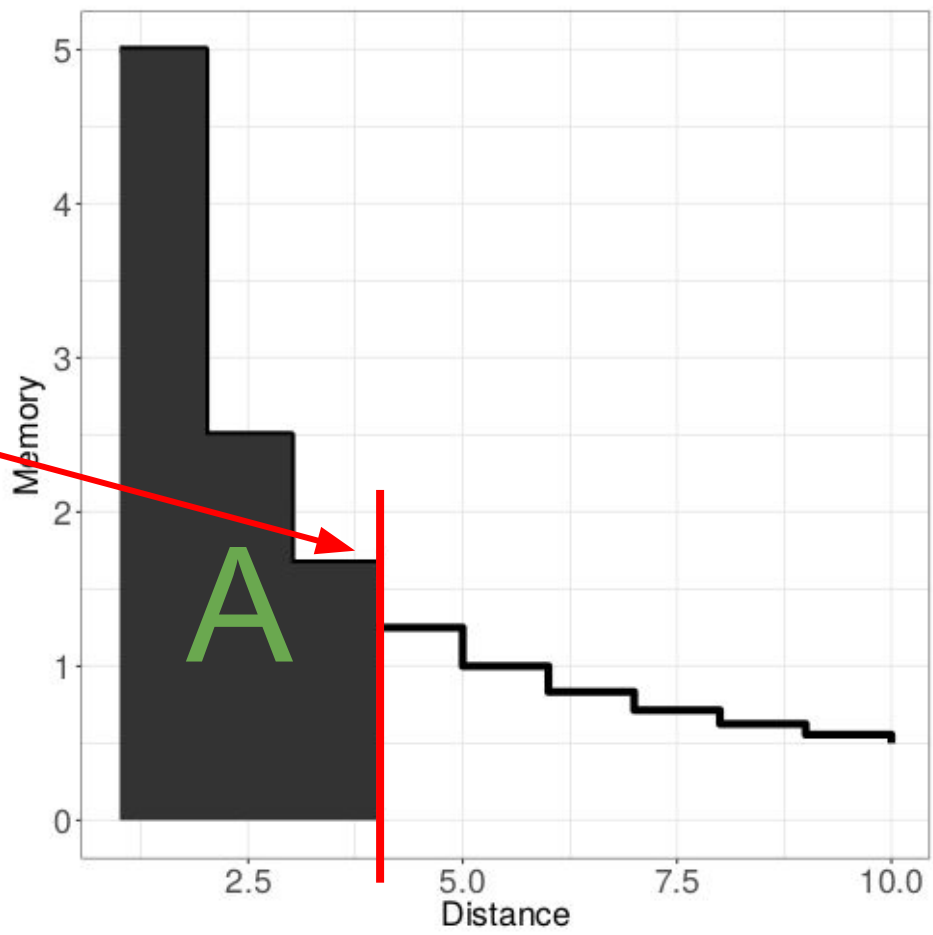
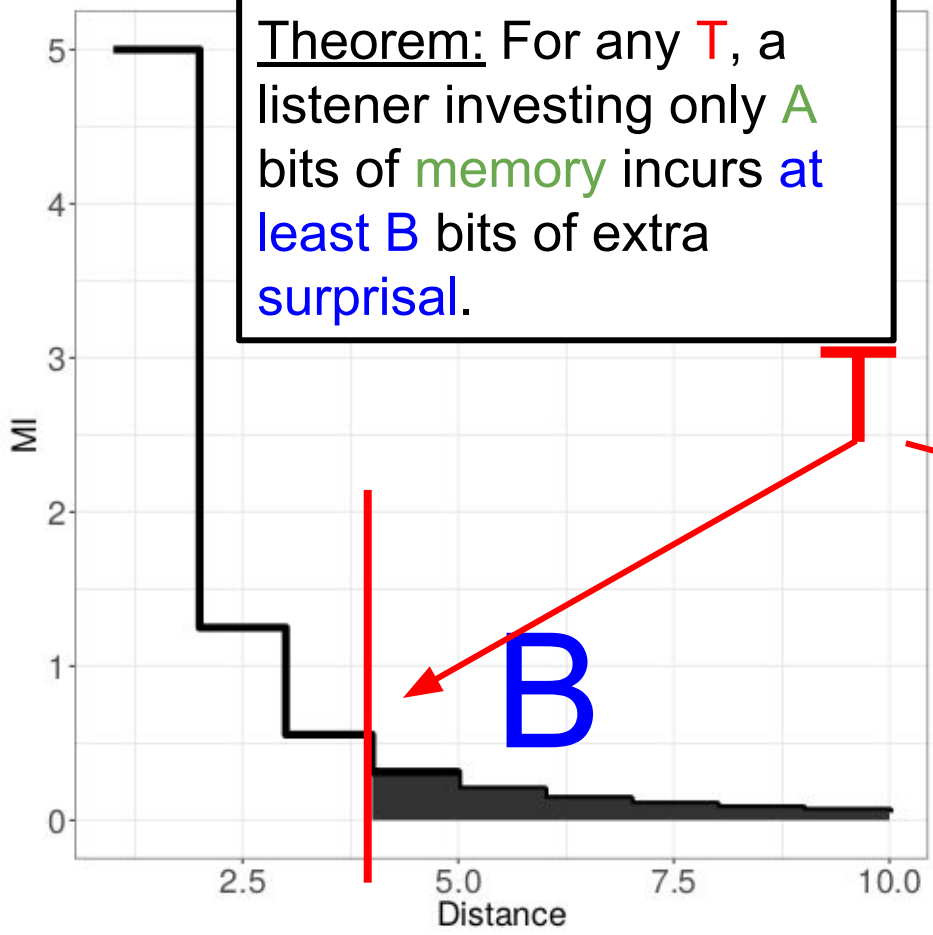
X_3

X_4

X_5

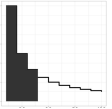
X_6

Theorem: For any T , a listener investing only A bits of memory incurs at least B bits of extra surprisal.



Proof

Assume that the listener's memory contains at most



bits



X_{-6}

X_{-5}

X_{-4}

X_{-3}

X_{-2}

X_{-1}

X_1

X_2

X_3

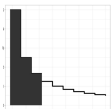
X_4

X_5

X_6

Proof

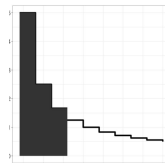
Assume that the listener's memory contains at most



bits



$$H[\text{man}] \leq$$



X_{-6}

X_{-5}

X_{-4}

X_{-3}

X_{-2}

X_{-1}

X_1

X_2

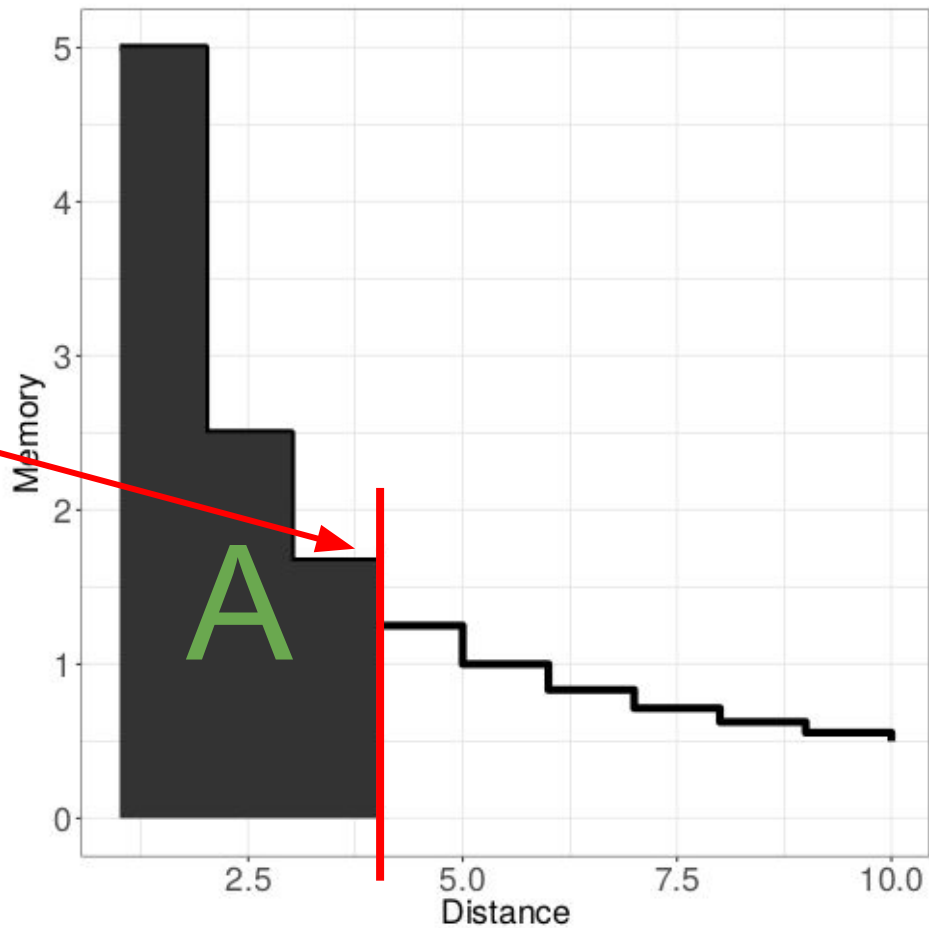
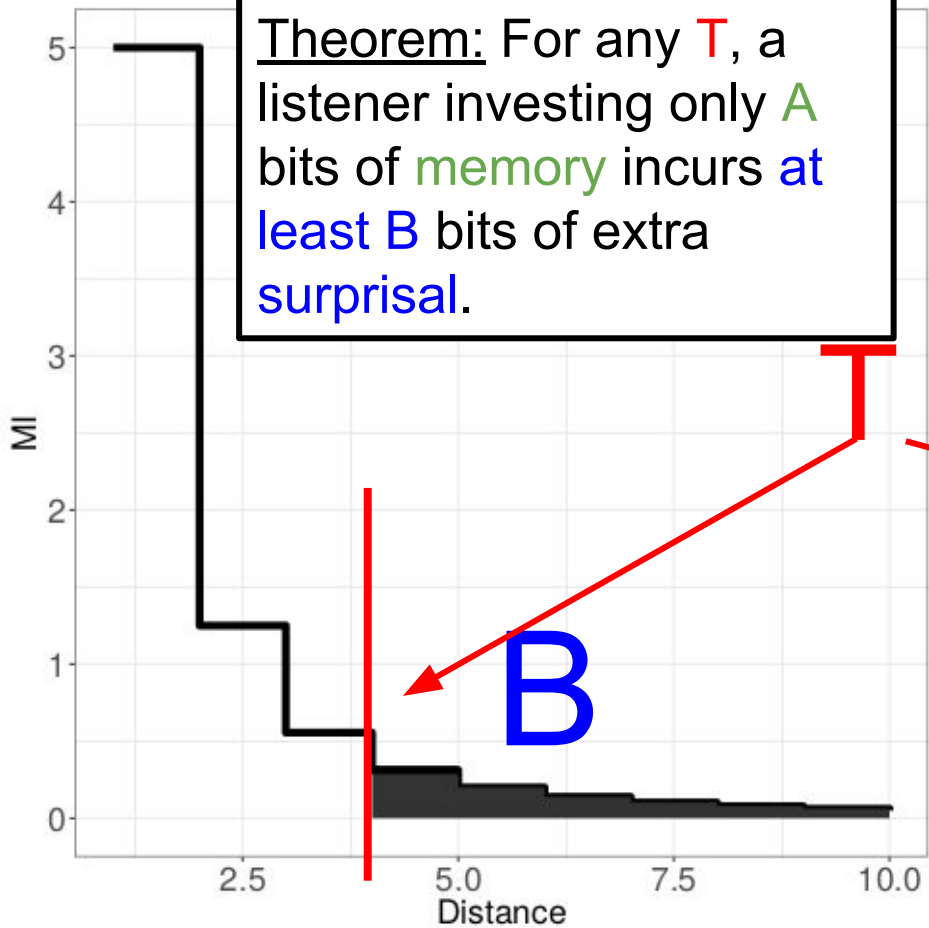
X_3

X_4

X_5

X_6

Theorem: For any T , a listener investing only A bits of memory incurs at least B bits of extra surprisal.



$$\text{Listener Surprisal} = H[X_1] - I[X_1, \text{Speaker}_0]$$

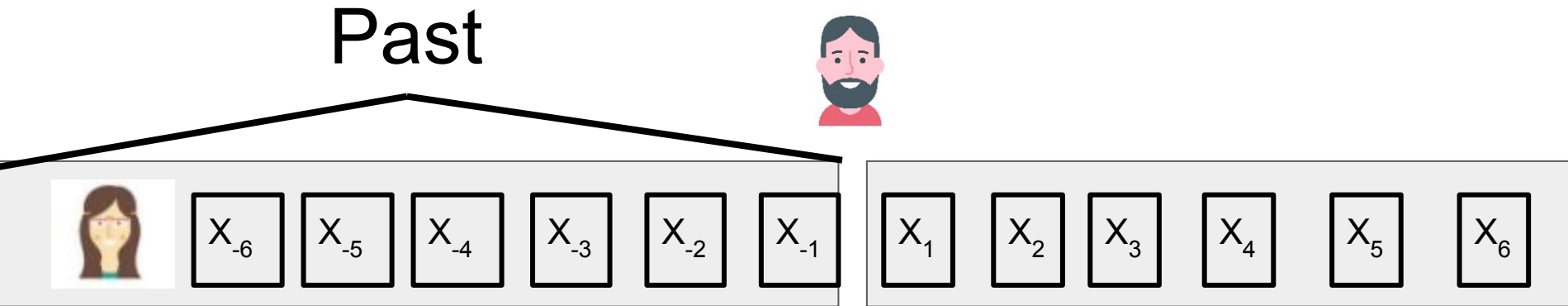


X_{-6} X_{-5} X_{-4} X_{-3} X_{-2} X_{-1}

X_1 X_2 X_3 X_4 X_5 X_6

$$\text{Listener Surprisal} = H[X_1] - I[X_1, \text{Person}_0]$$

$$\text{Optimal Surprisal} = H[X_1] - I[X_1, \text{Past}]$$



$$\text{Listener Surprisal} = H[X_1] - I[X_1, \text{Person}_0]$$

$$\text{Optimal Surprisal} = H[X_1] - I[X_1, \text{Past}]$$

Listener's **extra** surprisal is equal to

$$I[X_1, \text{Past}] - I[X_1, \text{Person}_0]$$

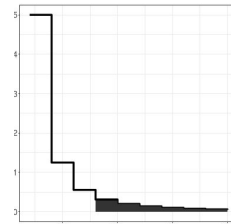
$$\text{Listener Surprisal} = H[X_1] - I[X_1, \text{Person}_0]$$

$$\text{Optimal Surprisal} = H[X_1] - I[X_1, \text{Past}]$$

Listener's **extra** surprisal is equal to

$$I[X_1, \text{Past}] - I[X_1, \text{Person}_0]$$

We want to lower-bound this by



Listener's extra surprisal is equal to

$$I[X_1, \text{Past}] - I[X_1, \text{Man} \text{ } 0]$$

Listener's extra surprisal is equal to

$$I[X_1, \text{Past}] - I[X_1, \text{Man}_0]$$

Bound this by averaging over a block of T words:

$$I[X_1, \text{Past}] - I[X_1, \text{Man}_0] \geq \frac{1}{T} (I[X_{1...T} | \text{Past}] - I[X_{1...T} | \text{Man}_0])$$

Listener's extra surprisal is equal to

$$I[X_1, \text{Past}] - I[X_1, \text{Past}, \text{Listener}]$$

Bound this by averaging over a block of T words:

$$I[X_1, \text{Past}] - I[X_1, \text{Past}, \text{Listener}] \geq \frac{1}{T} \left(I[X_1 \dots T | \text{Past}] - I[X_1 \dots T | \text{Past}, \text{Listener}] \right)$$

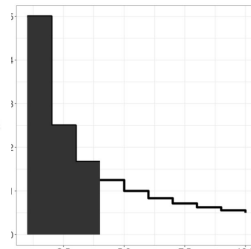
This is bounded by the listener's memory!

Listener's extra surprisal is equal to

$$I[X_1, \text{Past}] - I[X_1, \text{Person}_0]$$

Bound this by averaging over a block of T words:

$$I[X_1, \text{Past}] - I[X_1, \text{Person}_0] \geq \frac{1}{T} \left(I[X_{1...T} | \text{Past}] - \right)$$



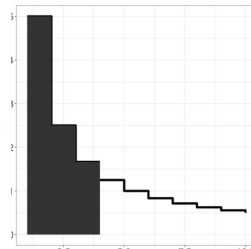
Lower-bound on
listener memory

Listener's extra surprisal is equal to

$$I[X_1, \text{Past}] - I[X_1, \text{Man} \text{ }_0]$$

Bound this by averaging over a block of T words:

$$I[X_1, \text{Past}] - I[X_1, \text{Man} \text{ }_0] \cong \frac{1}{T} \left(I[X_{1...T} | \text{Past}] - \right)$$



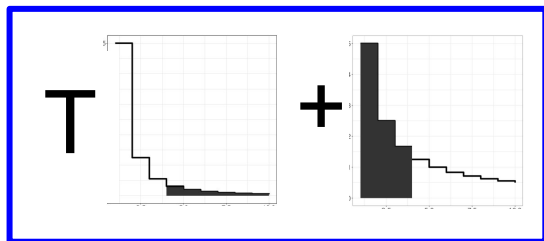
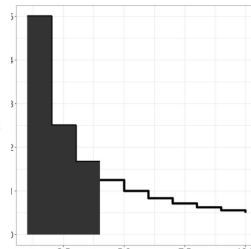
Can compute this explicitly

Listener's extra surprisal is equal to

$$I[X_1, \text{Past}] - I[X_1, \text{Person}_0]$$

Bound this by averaging over a block of T words:

$$I[X_1, \text{Past}] - I[X_1, \text{Person}_0] \geq \frac{1}{T} \left(I[X_{1...T} | \text{Past}] - \text{Person}_0 \right)$$

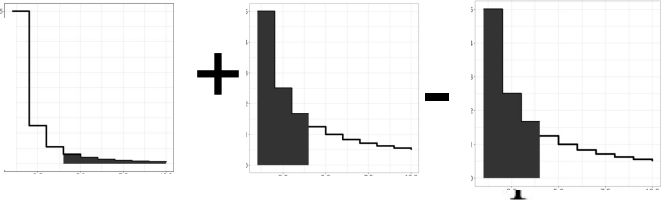


Can compute this explicitly

Listener's extra surprisal is equal to

$$I[X_1, \text{Past}] - I[X_1, \text{👤}_0]$$

Bound this by averaging over a block of T words:

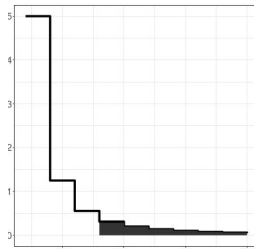
$$I[X_1, \text{Past}] - I[X_1, \text{👤}_0] \cong \frac{1}{T} \left(T \left[\text{Histogram 1} + \text{Histogram 2} - \text{Histogram 3} \right] \right)$$


Listener's extra surprisal is equal to

$$I[X_1, \text{Past}] - I[X_1, \text{👤}_0]$$

Bound this by averaging over a block of T words:

$$I[X_1, \text{Past}] - I[X_1, \text{👤}_0] \geq$$

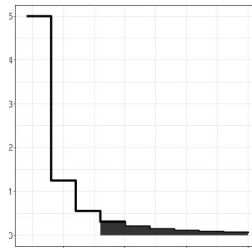


Listener's extra surprisal is equal to

$$I[X_1, \text{Past}] - I[X_1, \text{👤}_0]$$

Bound this by averaging over a block of T words:

$$I[X_1, \text{Past}] - I[X_1, \text{👤}_0] \geq$$



QED