

# Supplementary Information for: Crosslinguistic Word Orders Enable an Efficient Tradeoff between Memory and Surprisal

Michael Hahn, Judith Degen, Richard Futrell

2018

## 1 Formal Analysis and Proofs

In this section, we prove Theorem 1.

### 1.1 Mathematical Assumptions

We first make explicit how we formalize language processing for proving the theorem.

**Ingredient 1: Language as a Stationary Stochastic Process** We represent language as a stochastic process of words  $\dots w_{-2}w_{-1}w_0w_1w_2\dots$ , extending indefinitely both into the past and into the future. The symbols  $w_i$  belong to a common set, representing the words of the language.<sup>1</sup>

The assumption of infinite length is for mathematical convenience and does not affect the substance of our results: As we restrict our attention to the processing of individual sentences, which have finite length, we will actually not make use of long-range and infinite contexts.

We make the assumption that this process is *stationary*. Formally, this means that the conditional distribution  $P(w_t|w_{<t})$  does not depend on  $t$ , it only depends on the (semi-infinite) context sequence  $w_{<t}$ . Informally, this says that the process has no ‘internal clock’, and that the statistical rules of the language do not change at the timescale we are interested in. In reality, the statistical rules of language do change: They change as language changes over generations, and they also change between different situations – e.g., depending on the interlocutor at a given point in time. Given that we are interested in memory needs in the processing of *individual sentences*, at a timescale of seconds or minutes, stationarity seems to be a reasonable assumption to make.

**Ingredient 2: Postulates about Processing** The second ingredient consists of the three postulates described in the main paper. There are no further assumptions about the memory architecture and the nature of its computations.

### 1.2 Proof of the Theorem

We restate the theorem:

---

<sup>1</sup>Could also be phonemes, sentences, ..., any other kind of unit.

**Theorem 1.** Let  $T$  be any positive integer ( $T \in \{1, 2, 3, \dots\}$ ), and consider a listener using at most

$$\sum_{t=1}^T tI_t \quad (1)$$

bits of memory on average. Then this listener will incur surprisal at least

$$H[w_t|w_{<t}] + \sum_{t>T} I_t$$

on average.

*Proof.* The difference between the listener's surprisal and optimal surprisal is  $H[w_t|m_t] - H[w_t|w_{<t}]$ .<sup>2</sup> By the assumption of stationarity, we can, for any positive integer  $T$ , rewrite this expression as

$$H[w_t|m_t] - H[w_t|w_{<t}] = \frac{1}{T} \sum_{t'=1}^T (H[w_{t'}|m_{t'}] - H[w_{t'}|w_{<t'}]) \quad (2)$$

Because  $m_t$  is determined by  $(w_{1\dots t-1}, m_1)$ :

$$m_t = M(m_{t-1}, w_{t-1}) = M(M(m_{t-2}, w_{t-2}), w_{t-1}) = M(M(M(m_{t-3}, w_{t-3}), w_{t-2}), w_{t-1}) = \dots \quad (3)$$

the Data Processing inequality entails the following inequality for every positive integer  $t$ :

$$H[w_t|m_t] \geq H[w_t|w_{1\dots t-1}, m_1] \quad (4)$$

Plugging this inequality into Equation 2 above:

$$H[w_t|m_t] - H[w_t|w_{<t}] \geq \frac{1}{T} \sum_{t=1}^T (H[w_t|w_{1\dots t-1}, m_1] - H[w_t|w_{1\dots t-1}, w_{\leq 0}]) \quad (5)$$

$$= \frac{1}{T} (H[w_{1\dots T}|m_1] - H[w_{1\dots T}|w_{\leq 0}]) \quad (6)$$

$$= \frac{1}{T} (I[w_{1\dots T}, w_{\leq 0}] - I[w_{1\dots T}, m_1]) \quad (7)$$

The first term  $I[w_{1\dots T}, w_{\leq 0}]$  can be rewritten in terms of  $I_t$ :

$$I[w_{1\dots T}, w_{\leq 0}] = \sum_{i=1}^T \sum_{j=-1}^{-\infty} I[w_i, w_j | w_{j+1} \dots w_{i-1}] = \sum_{t=1}^T tI_t + T \sum_{t>T} I_t \quad (8)$$

Therefore

$$H[w_t|m_t] - H[w_t|w_{<t}] \geq \frac{1}{T} \left( \sum_{t=1}^T tI_t + T \sum_{t>T} I_t - I[w_{1\dots T}, m_1] \right)$$

$I[w_{1\dots T}|m_1]$  is at most  $H[m_1]$ , which is at most  $\sum_{t=1}^T tI_t$  by assumption. Thus, the expression above is bounded by

$$\begin{aligned} H[w_t|m_t] - H[w_t|w_{<t}] &\geq \frac{1}{T} \left( \sum_{t=1}^T tI_t + T \sum_{t>T} I_t - \sum_{t=1}^T tI_t \right) \\ &= \sum_{t>T} I_t \end{aligned}$$

Rearranging shows that the listener's surprisal is at least  $H[w_t|m_t] \geq H[w_t|w_{<t}] + \sum_{t>T} I_t$ , as claimed.  $\square$

<sup>2</sup>A listener whose predictions are not optimal given  $m_t$  can only incur even higher surprisal.

### 1.3 For nondeterministic encoding functions

We have been assuming that  $m_t$  is a deterministic function of  $x_t$  and  $m_{t-1}$ . Here, we show that this assumption can be relaxed to stochastic encoding functions.

We relax condition (X) to the following requirement, for all values of  $m_1, (w_t)_{t \in \mathbb{Z}}, m_0$ :

$$p(m_1 | (w_t)_{t \in \mathbb{Z}}, m_0) = p(m_1 | m_0, w_1) \quad (9)$$

This says that  $m_1$  contains no information about the utterances beyond what is contained in  $m_0$  and  $w_1$ .

The one place in the proof where (X) plays a role is the proof of the inequality:

$$H[w_t | m_t] \geq H[w_t | w_{1..t-1}, m_1] \quad (10)$$

We show that this inequality still holds under the relaxed condition (9):

*Proof.* By Bayes' Theorem

$$\begin{aligned} p(w_t | m_0, m_1, w_{0..t-1}) &= \frac{p(m_1 | m_0, w_{0..t})}{p(m_1 | m_0, w_{0..t-1})} \cdot p(w_t | m_0, w_{0..t-1}) \\ &= \frac{p(m_1 | m_0, w_0)}{p(m_1 | m_0, w_0)} \cdot p(w_t | m_0, w_{0..t-1}) \\ &= p(w_t | m_0, w_{0..t-1}) \end{aligned}$$

where the second equation follows from (9). So we have a Markov chain

$$(w_t) \rightarrow (m_0, w_{0..t-1}) \rightarrow (m_1, w_{1..t-1}) \quad (11)$$

Thus, by the Data Processing Inequality,

$$H[w_t | w_{1..t-1}, m_1] \geq H[w_t | w_{0..t-1}, m_0] \quad (12)$$

Finally, iteratively applying this reasoning, we conclude:

$$H[w_t | m_t] \geq H[w_t | w_{t-1}, m_{t-1}] \geq H[w_t | w_{t-2}, m_{t-2}] \geq \dots \geq H[w_t | w_{1..t-1}, m_1]$$

□

### 1.4 Locality in a model with Memory Retrieval

Here we show that our information-theoretic analysis is compatible with models placing the main bottleneck in the difficulty of retrieval (McElree, 2000; Lewis and Vasishth, 2005; Nicenboim and Vasishth, 2018; Vasishth et al., 2019). We extend our model of memory in incremental prediction to capture key aspects of the models described by Lewis and Vasishth (2005); Nicenboim and Vasishth (2018); Vasishth et al. (2019).

The ACT-R model of Lewis and Vasishth (2005) assumes a small working memory consisting of *buffers* and a *control state*, which together hold a small and fixed number of individual *chunks*. It also assumes a large short-term memory that contains an unbounded number of chunks. This large memory store is accessed via *cue-based retrieval*: a query is constructed based on the current state of the buffers and the control state; a chunk that matches this query is then selected from the memory storage and placed into one of the buffers.

**Formal Model** We extend our information-theoretic analysis by considering a model that maintains both a small working memory  $m_t$  – corresponding to the buffers and the control state – and an unlimited short-term memory  $s_t$ . Predictions are made based on working memory  $m_t$ , incurring surprisal  $H[w_t|m_t]$ . When processing a word  $x_t$ , there is some amount of communication between  $m_t$  and  $s_t$ , corresponding to retrieval operations. We model this using a variable  $r_t$  representing the information that is retrieved from  $s_t$ . In our formalization,  $r_t$  reflects the totality of all retrieval operations that are made during the processing of  $x_{t-1}$ ; they happen after  $x_{t-1}$  has been observed but before  $x_t$  has.

The working memory state is determined not just by the input  $x_t$  and the previous working memory state  $m_{t-1}$ , but also by the retrieved information:

$$m_t = f(x_t, m_{t-1}, r_t) \quad (13)$$

The retrieval operation is jointly determined by working memory, short-term memory, and the previous word:

$$r_t = g(x_{t-1}, m_{t-1}, s_{t-1}) \quad (14)$$

Finally, the short-term memory can incorporate any – possibly all – information from the last word and the working memory:

$$s_t = h(x_{t-1}, m_{t-1}, s_{t-1}) \quad (15)$$

While  $s_t$  is unconstrained, there are constraints on the capacity of working memory  $H[m_t]$  and the amount of retrieved information  $H[r_t]$ . Placing a bound on  $H[m_t]$  reflects the fact that the buffers can only hold a small and fixed number of chunks (Lewis and Vasishth, 2005).

**Cost of Retrieval** In the model of Lewis and Vasishth (2005), the time it takes to process a word is determined primarily by the time spent retrieving chunks, which is determined by the number of retrieval operations and the time it takes to complete each retrieval operation. If the information content of each chunk is bounded, then a bound on  $H[r_t]$  corresponds to a bound on the number of retrieval operations.

In the model of Lewis and Vasishth (2005), a retrieval operation takes longer if more chunks are similar to the retrieval cue, whereas, in the direct-access model (McElree, 2000; Nicenboim and Vasishth, 2018; Vasishth et al., 2019), retrieval operations take a constant amount of time. There is no direct counterpart to differences in retrieval times and similarity-based inhibition as in the activation-based model in our formalization. Our formalization thus more closely matches the direct-access model, though it might be possible to incorporate aspects of the activation-based model in our formalization.

**Role of Surprisal** The ACT-R model of Lewis and Vasishth (2005) does not have an explicit surprisal cost. Instead, surprisal effects are interpreted as arising because, in less constraining contexts, the parser is more likely to make decisions that then turn out to be incorrect, leading to additional correcting steps. We view this as an algorithmic-level implementation of a surprisal cost  $H[x_t|m_{t-1}]$ : If the word  $x_t$  is unexpected given the current state of the working memory – i.e., buffers and control states – then their current state must provide insufficient information to constrain the actual syntactic state of the sentence, meaning that the parsing steps made to integrate  $x_t$  are likely to include more backtracking and correction steps. Thus, we argue that cue-based retrieval models predict that the surprisal  $-\log P(x_t|m_{t-1})$  will be part of the cost of processing word  $x_t$ .

**Theoretical Result** We now show an extension of our theoretical result in the setting of the retrieval-based model described above.

**Theorem 2.** Let  $0 < S \leq T$  be positive integers such that the average working memory cost  $H[m_t]$  is bounded as

$$H[m_t] \leq \sum_{t=1}^T tI_t \quad (16)$$

and the average amount of retrieved information is bounded as

$$H[r_t] \leq \sum_{t=T+1}^S I_t \quad (17)$$

Then the surprisal cost is lower-bounded as

$$H[w_t|m_t] \geq H[w_t|x_{<t}] + \sum_{t>S} I_t \quad (18)$$

*Proof.* The proof is a generalization of the proof above. For any positive integer  $t$ ,  $m_t$  is determined by  $w_{1...t}, m_0, r_0, \dots, r_t$ . Therefore, the Data Processing Inequality entails:

$$H[w_t|m_t] \geq H[w_t|w_{1...t}, m_0, r_0, \dots, r_t] \quad (19)$$

As in (5), this leads to

$$H[w_t|m_t] - H[w_t|w_{<t}] \geq \frac{1}{T} \sum_{t=1}^T (H[w_t|w_{1...t}, m_0, r_0, \dots, r_t] - H[w_t|w_{1...t-1}, w_{\leq 0}]) \quad (20)$$

$$\geq \frac{1}{T} (H[w_{1...T}|m_0, r_0, \dots, r_T] - H[w_{1...T}|w_{\leq 0}]) \quad (21)$$

$$= \frac{1}{T} (I[w_{1...T}, w_{\leq 0}] - I[w_{1...T}, (m_0, r_0, \dots, r_T)]) \quad (22)$$

Now, using the calculation from (8), this can be rewritten as:

$$\begin{aligned} H[w_t|m_t] - H[w_t|w_{<t}] &= \frac{1}{T} \left( \sum_{t=1}^T tI_t + T \sum_{t>T} I_t - I[X_1 \dots X_T, (M_0, R_1, \dots, R_T)] \right) \\ &= \frac{1}{T} \left( \sum_{t=1}^T tI_t + T \sum_{t>T} I_t - I[X_{1...T}, M_0] - \sum_{t=1}^T I[X_{1...T}, R_t|M_0, r_{1...t-1}] \right) \end{aligned}$$

Due to the inequalities  $I[X_{1...T}, M_0] \leq H[M_0]$  and  $I[X_{1...T}, R_t|M_0, r_{1...t-1}] \leq H[R_t]$ , this can be bounded as

$$H[w_t|m_t] - H[w_t|w_{<t}] \geq \frac{1}{T} \left( \sum_{t=1}^T tI_t + T \sum_{t>T} I_t - H[M_0] - \sum_{t=1}^T H[R_t] \right) \quad (23)$$

$$(24)$$

Finally, this reduces as

$$H[w_t|m_t] - H[w_t|w_{<t}] \geq \frac{1}{T} (T \sum_{t>T} I_t - T \cdot H[R_t]) \quad (25)$$

$$= \sum_{t>T} I_t - H[R_t] \quad (26)$$

$$\geq \sum_{t>T} I_t - \sum_{t=T+1}^S I_t \quad (27)$$

$$= \sum_{t>S} I_t \quad (28)$$

□

**Information Locality** We now show that this result predicts information locality provided that retrieving information is more expensive than keeping the same amount of information in working memory. For this, we formalize the problem of finding an optimal memory strategy as a multi-objective optimization, aiming to minimize

$$\lambda_1 H[m_t] + \lambda_2 H[r_t] \quad (29)$$

to achieve a given surprisal level, for some setting of  $\lambda_1, \lambda_2 > 0$  describing the relative cost of storage and retrieval. What is the optimal division of labor between keeping information in working memory and recovering it through retrieval? The problem

$$\min_T \lambda_1 \sum_{t=1}^T t I_t + \lambda_2 \sum_{t=T+1}^S I_t \quad (30)$$

has solution  $T \approx \frac{\lambda_2}{\lambda_1}$ . This means that, as long as retrievals are more expensive than keeping the same amount of information in working memory (i.e.,  $\lambda_2 > \lambda_1$ ), the optimal strategy stores information from the last  $T > 1$  words in working memory. Due to the factor  $t$  inside  $\sum_{t=1}^T t I_t$ , the bound (30) will be reduced when  $I_t$  decays faster, i.e., there is strong information locality.

The assumption that retrieving information is more difficult than storing it is reasonable for cue-based retrieval models, as retrieval suffers from similarity-based interference effects due to the unstructured nature of the storage (Lewis and Vasishth, 2005). A model that maintains no information in its working memory, i.e.  $H[m_t] = 0$ , would correspond to a cue-based retrieval model that stores nothing in its buffers and control states, and relies entirely on retrieval to access past information. Given the nature of representations assumed in models (Lewis and Vasishth, 2005), such a model would seem to be severely restricted in its ability to parse language.

## 1.5 Results for Language Production

Here we show results linking memory and locality in production. We show that results similar to our main theorem hold for the tradeoff between a speaker’s memory and the accuracy with which they match the distribution of the language.

**Speaker aims to match language distribution** First, we consider a setting in which a speaker produces sentences with bounded memory, and analyze the deviation of the produced distribution from the actual distribution of the language.

We consider a speaker who maintains memory representations and incrementally produces based on these representations:

$$p_{\text{speaker}}(x_t|X_{<t}) = p(x_t|m_t) \quad (31)$$

We show a tradeoff between the memory capacity  $H[m_t]$  and the KL-divergence between the actual language statistics and the speaker's production distribution:

$$D_{KL}(P_{\text{language}}||P_{\text{produced}}) := \mathbb{E}_{X_{<t}} \sum_{x_t} p(x_t|X_{<t}) \log \frac{p(x_t|X_{<t})}{p_{\text{speaker}}(x_t|X_{<t})} \quad (32)$$

**Theorem 3.** *If a speaker maintains memory*

$$H[m_t] \leq \sum_{i=1}^T t I_t \quad (33)$$

then

$$D_{KL}(P_{\text{language}}||P_{\text{produced}}) \geq \sum_{t=T+1}^{\infty} I_t \quad (34)$$

While this bound only considers the production of a single word, it immediately entails a bound on the production accuracy for sequences:

$$D_{KL}(P_{\text{language}}(X_1 \dots X_t|X_{\leq 0})||P_{\text{produced}}(X_1 \dots X_t|X_{\leq 0})) = t \cdot D_{KL}(P_{\text{language}}(X_1|X_{\leq 0})||P_{\text{produced}}(X_1|X_{\leq 0})) \quad (35)$$

*Proof.* First note

$$D_{KL}(P_{\text{language}}||P_{\text{produced}}) = \mathbb{E}_{X_{<t}} \sum_{x_t} p(x_t|X_{<t}) \log \frac{p(x_t|X_{<t})}{p_{\text{speaker}}(x_t|X_{<t})} \quad (36)$$

$$= \mathbb{E}_{X_{<t}} \sum_{x_t} p(x_t|X_{<t}) \log \frac{p(x_t|X_{<t})}{p(x_t|M(X_{<t}))} \quad (37)$$

$$= \mathbb{E}_{X_{<t}} \sum_{x_t} p(x_t|X_{<t}) \log p(x_t|X_{<t}) - \mathbb{E}_{X_{<t}} \sum_{x_t} p(x_t|X_{<t}) \log p(x_t|M(X_{<t})) \quad (38)$$

$$= \mathbb{E}_{X_{<t}} \sum_{x_t} p(x_t|X_{<t}) \log p(x_t|X_{<t}) + S_M(x_t|X_{<t}) \quad (39)$$

In the last line, the first term is a constant independent of  $M$ .

Then the proof for the listener case transfers to this setting.  $\square$

TODO limitations of this as a model of production

**Speaker aims to match, conditional on goal** The first setting does not account for the fact that language is produced aiming for some communicative goal. We therefore now assume that the speaker has a communicative goal  $G$  in mind. This goal  $G$  stays constant during production process for a sentence, and we count how much memory is needed in addition to the goal  $G$ . We assume that there is a distribution of sentences expressing goals  $G$ :

$$P(\text{sentence}|G) \quad (40)$$

and assume that the speaker aims to match this distribution

$$\mathbb{E}_G[D_{KL}((language|G)||(\textit{produced}|G)))] \quad (41)$$

We can analyze this model by adding conditioning w.r.t.  $G$  throughout the analysis of the previous case. Specifically, we need

$$I_t^G := I[X_t, X_0 | X_1, \dots, X_{t-1}, G] \quad (42)$$

Take  $I_t$  conditioned on  $G$ : only count statistical dependencies to the degree that they are not redundant with the goal

**Theorem 4.** *If a speaker maintains memory*

$$H[m_t] \leq \sum_{i=1}^T t I_t^G \quad (43)$$

then

$$\mathbb{E}_G D_{KL}(P_{language}(\cdot|G) || P_{produced}(\cdot|G)) \geq \sum_{t=T+1}^{\infty} I_t^G \quad (44)$$

*Proof.* This is entirely analogous to the previous proof.  $\square$

are there conditions under which this is close to  $I_t$ ?

**Pragmatic Speaker** would need an assumption on the density of goals in the space of sequences.

Note

$$D_{KL}(P_{language}(x_{1...t}) || P_{produced}(x_{1...t})) := \sum_{x_{1...t}} p(x_t | X_{1...t}) \log \frac{p(x_t | X_{1...t})}{p_{speaker}(x_t | X_{1...t})} \geq t D_{KL}(x_t || \dots) \quad (45)$$

$$H[G | Produced] - H[G | Language]$$

## 2 Proof of Left-Right Invariance

Here we show that the bound provided by our theorem is invariant under reversal of the process. That is: Given a process  $(X_t)_{t \in \mathbb{Z}}$ , we define its reverse process  $(Y_t)_{t \in \mathbb{Z}}$  by  $Y_t := X_t$ . We claim that the theorem provides the same bounds for the memory-surprisal tradeoff curves. To prove this, we note:

$$I[X_t, X_0 | X_{1...t-1}] = I[Y_{-t}, Y_0 | Y_{1-t...-1}] = I[Y_0, Y_t | Y_{1...t-1}] = I[Y_t, Y_0 | Y_{1...t-1}] \quad (46)$$

The first step follows from the definition of  $Y$ . The second step follows from the fact that  $X_t$ , and thus also  $Y_t$ , is stationary, and thus adding  $t$  to each index in the expression does not change the resulting value. The third step uses the fact that mutual information is symmetric.



### 3 Example where window model is not optimal

Here we provide an example of a stochastic process where a window-based memory encoding is not optimal, but the bound provided by our theorem still holds.

Let  $k$  be some positive integer. Consider a process  $x_{t+1} = (v_{t+1}, w_{t+1}, y_{t+1}, z_{t+1})$  where

1. The first two components consist of fresh random bits. Formally,  $v_{t+1}$  is an independent draw from  $Bernoulli(0.5)$ , independent from all preceding observations  $x_{\leq t}$ . Second, let  $w_{t+1}$  consist of  $2k$  many such independent random bits (so that  $H[w_{t+1}] = 2k$ )
2. The third component *deterministically* copies the first bit from  $2k$  steps earlier. Formally,  $y_{t+1}$  is equal to the first component of  $x_{t-2k+1}$
3. The fourth component *stochastically* copies the second part (consisting of  $2k$  random bits) from one step earlier. Formally, each component  $z_{t+1}^{(i)}$  is determined as follows: First take a sample  $u_{t+1}^{(i)}$  from  $Bernoulli(\frac{1}{4k})$ , independent from all preceding observations. If  $u_{t+1}^{(i)} = 1$ , set  $z_{t+1}^{(i)}$  to be equal to the second component of  $w_t^{(i)}$ . Otherwise, let  $z_{t+1}^{(i)}$  be a fresh draw from  $Bernoulli(0.5)$ .

Predicting observations optimally requires taking into account observations from the  $2k$  last time steps.

We show that, when approximately predicting with low memory capacities, a window-based approach does *not* in general achieve an optimal memory-surprisal tradeoff.

Consider a model that predicts  $x_{t+1}$  from only the last observation  $x_t$ , i.e., uses a window of length one. The only relevant piece of information in this past observation is  $w_t$ , which stochastically influences  $z_{t+1}$ . Storing this costs  $2k$  bit of memory as  $w_t$  consists of  $2k$  draws from  $Bernoulli(0.5)$ . How much does it reduce the surprisal of  $x_{t+1}$ ? Due to the stochastic nature of  $z_{t+1}$ , it reduces the surprisal only by about  $I[x_{t+1}, w_t] = I[z_{t+1}, w_t] < 2k \cdot \frac{1}{2k} = 1$ , i.e., surprisal reduction is strictly less than one bit.<sup>3</sup>

We show that there is an alternative model that strictly improves on this window-based model: Consider a memory encoding model that encodes each of  $v_{t-2k+1}, \dots, v_t$ , which costs  $2k$  bits of memory – as the window-based model did. Since  $y_{t+1} = v_{t-2k+1}$ , this model achieves a surprisal reduction of  $H[v_{t-2k+1}] = 1$  bit, strictly more than the window-based model.

This result does not contradict our theorem because the theorem only provides *bounds* across models, which are not necessarily achieved by a given window-based model. In fact, for the process described here, no memory encoding function  $M$  can exactly achieve the theoretical bound described by the theorem.

### 4 Corpus Size per Language

Language	Training	Held-Out	Language	Training	Held-Out
Afrikaans	1,315	194	Indonesian	4,477	559
Amharic	974	100	Italian	17,427	1,070
Arabic	21,864	2,895	Japanese	7,164	511
Armenian	514	50	Kazakh	947	100

<sup>3</sup>We can evaluate  $I[z_{t+1}, w_t]$  as follows. Set  $l = k/4$ . Write  $z, w$  for any of the  $2k$  components of  $z_{t+1}, w_t$ , respectively. First, calculate  $p(z = 1|w = 1) = 1/l + (1 - 1/l)\frac{1}{2} = 1/(2l) + 1/2 = \frac{1+l}{2l}$  and  $p(z = 0|w = 1) = (1 - 1/l)\frac{1}{2} = 1/2 - 1/2l = \frac{l-1}{2l}$ . Then  $I[Z, W] = D_{KL}(p(z|w = 1)||p(z)) = \frac{1+l}{2l} \log \frac{1+l}{1/2} + \frac{l-1}{2l} \log \frac{l-1}{1/2} = \frac{1+l}{2l} \log \frac{1+l}{l} + \frac{l-1}{2l} \log \frac{l-1}{l} \leq \frac{1+l}{l} \log \frac{1+l}{l} = (1 + 1/l) \log(1 + 1/l) \leq (1 + 1/l)(1/l) = 1/l + 1/l^2 < 2/l = \frac{1}{2k}$ .

Bambara	926	100	Korean	27,410	3,016
Basque	5,396	1,798	Kurmanji	634	100
Breton	788	100	Latvian	4,124	989
Bulgarian	8,907	1,115	Maltese	1,123	433
Buryat	808	100	Naija	848	100
Cantonese	550	100	North Sami	2,257	865
Catalan	13,123	1,709	Norwegian	29,870	4,639
Chinese	3,997	500	Persian	4,798	599
Croatian	7,689	600	Polish	6,100	1,027
Czech	102,993	11,311	Portuguese	17,995	1,770
Danish	4,383	564	Romanian	8,664	752
Dutch	18,310	1,518	Russian	52,664	7,163
English	17,062	3,070	Serbian	2,935	465
Erzya	1,450	100	Slovak	8,483	1,060
Estonian	6,959	855	Slovenian	7,532	1,817
Faroese	1,108	100	Spanish	28,492	3,054
Finnish	27,198	3,239	Swedish	7,041	1,416
French	32,347	3,232	Thai	900	100
German	13,814	799	Turkish	3,685	975
Greek	1,662	403	Ukrainian	4,506	577
Hebrew	5,241	484	Urdu	4,043	552
Hindi	13,304	1,659	Uyghur	1,656	900
Hungarian	910	441	Vietnamese	1,400	800

Table 2: Languages, with the number of training and held-out sentences available.

{tab:corpora

## 5 Samples Drawn per Language

Language	Base.	Real	Language	Base.	Real
Afrikaans	13	10	Indonesian	11	11
Amharic	137	10	Italian	10	10
Arabic	11	10	Japanese	25	15
Armenian	140	76	Kazakh	11	10
Bambara	25	29	Korean	11	10
Basque	15	10	Kurmanji	338	61
Breton	35	14	Latvian	308	178
Bulgarian	14	10	Maltese	30	24
Buryat	26	18	Naija	214	10
Cantonese	306	32	North Sami	335	194
Catalan	11	10	Norwegian	12	10
Chinese	21	10	Persian	25	12
Croatian	30	17	Polish	309	35
Czech	18	10	Portuguese	15	55

Danish	33	17	Romanian	10	10
Dutch	27	10	Russian	20	10
English	13	11	Serbian	26	11
Erzya	846	167	Slovak	303	27
Estonian	347	101	Slovenian	297	80
Faroese	27	13	Spanish	14	10
Finnish	83	16	Swedish	31	14
French	14	11	Thai	45	19
German	19	13	Turkish	13	10
Greek	16	10	Ukrainian	28	18
Hebrew	11	10	Urdu	17	10
Hindi	11	10	Uyghur	326	175
Hungarian	220	109	Vietnamese	303	12

Figure 1: Samples drawn per language according to the precision-dependent stopping criterion.

{tab:samples}

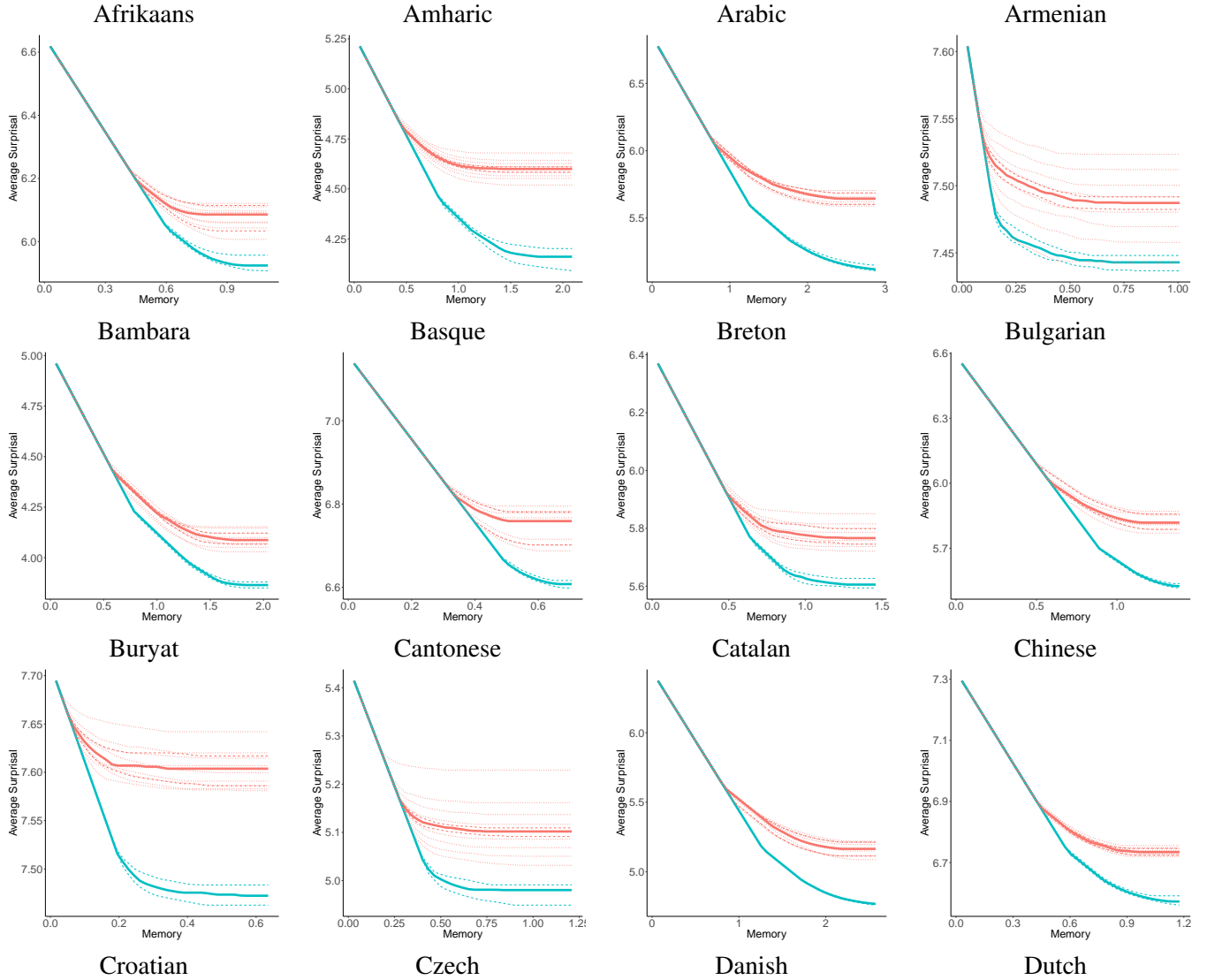
Language	Mean	Lower	Upper	Language	Mean	Lower	Upper
Afrikaans	1.0	1.0	1.0	Indonesian	1.0	1.0	1.0
Amharic	1.0	1.0	1.0	Italian	1.0	1.0	1.0
Arabic	1.0	1.0	1.0	Japanese	1.0	1.0	1.0
Armenian	0.92	0.87	0.97	Kazakh	1.0	1.0	1.0
Bambara	1.0	1.0	1.0	Korean	1.0	1.0	1.0
Basque	1.0	1.0	1.0	Kurmanji	0.93	0.88	0.98
Breton	1.0	1.0	1.0	Latvian	0.49	0.4	0.57
Bulgarian	1.0	1.0	1.0	Maltese	1.0	1.0	1.0
Buryat	1.0	1.0	1.0	Naija	1.0	0.99	1.0
Cantonese	0.96	0.86	1.0	North Sami	0.37	0.3	0.44
Catalan	1.0	1.0	1.0	Norwegian	1.0	1.0	1.0
Chinese	1.0	1.0	1.0	Persian	1.0	1.0	1.0
Croatian	1.0	1.0	1.0	Polish	0.1	0.04	0.17
Czech	1.0	1.0	1.0	Portuguese	1.0	1.0	1.0
Danish	1.0	1.0	1.0	Romanian	1.0	1.0	1.0
Dutch	1.0	1.0	1.0	Russian	1.0	1.0	1.0
English	1.0	1.0	1.0	Serbian	1.0	1.0	1.0
Erzya	0.99	0.98	1.0	Slovak	0.07	0.03	0.12
Estonian	0.8	0.72	0.86	Slovenian	0.82	0.77	0.88
Faroese	1.0	1.0	1.0	Spanish	1.0	1.0	1.0
Finnish	1.0	1.0	1.0	Swedish	1.0	1.0	1.0
French	1.0	1.0	1.0	Thai	1.0	1.0	1.0
German	1.0	0.91	1.0	Turkish	1.0	1.0	1.0
Greek	1.0	1.0	1.0	Ukrainian	1.0	1.0	1.0
Hebrew	1.0	1.0	1.0	Urdu	1.0	1.0	1.0
Hindi	1.0	1.0	1.0	Uyghur	0.65	0.57	0.73
Hungarian	0.87	0.8	0.93	Vietnamese	1.0	0.98	1.0

Figure 2: Bootstrapped estimates for  $G$ .

{tab:boot-g}

## 6 Detailed Results per Language

### 6.1 Median Surprisal per Memory Budget



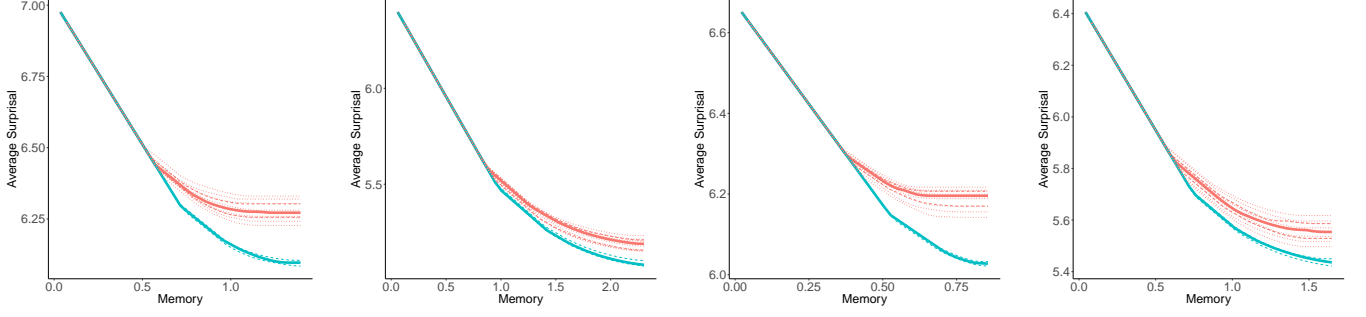
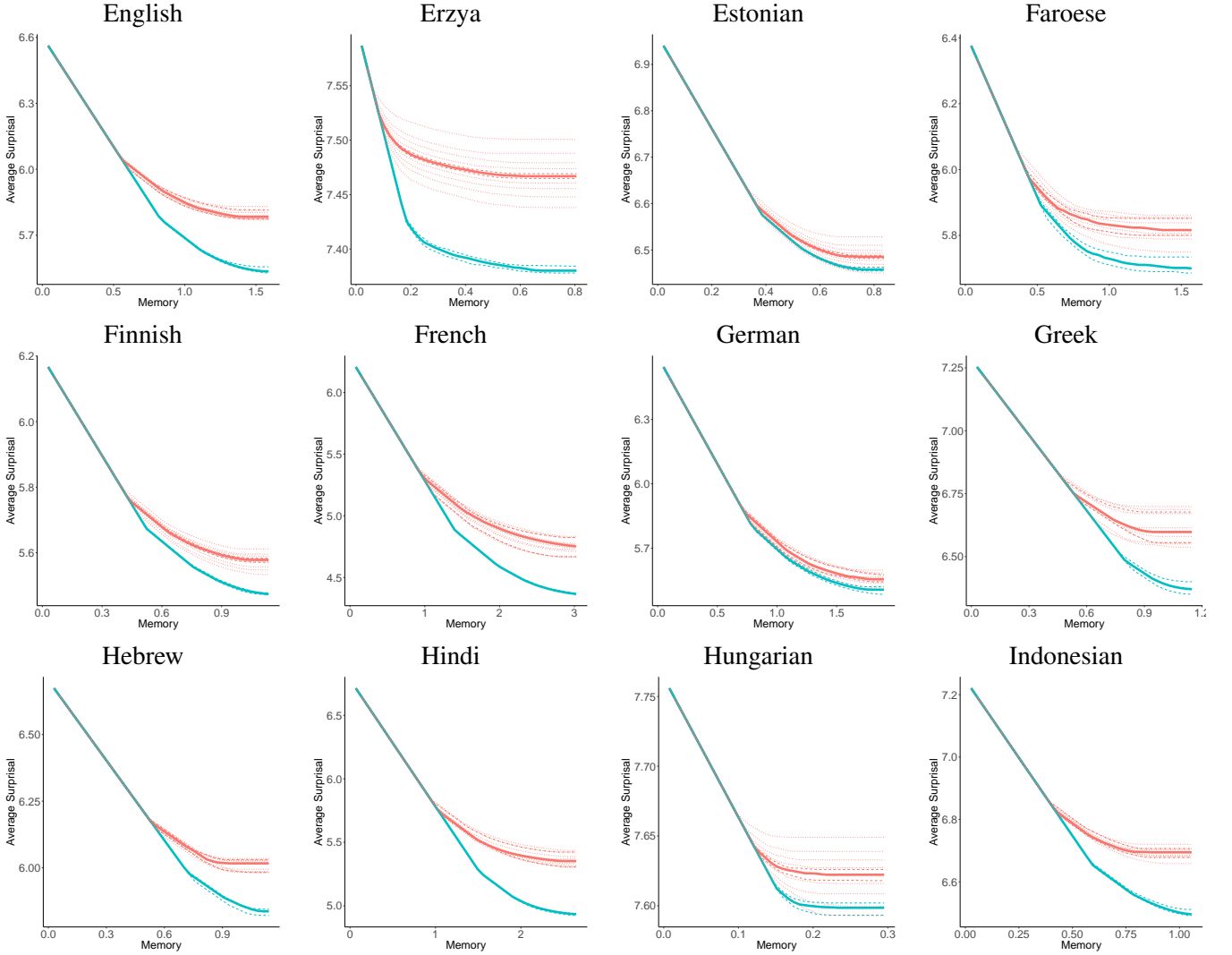


Figure 3: Medians: For each memory budget, we provide the median surprisal for real and random languages. Solid lines indicate sample medians, dashed lines indicate 95 % confidence intervals for the population median, dotted lines indicate empirical quantiles (10%, 20%, ..., 80%, 90%). Green: Random baselines; blue: real language; red: maximum-likelihood grammars fit to real orderings.

{tab:medians



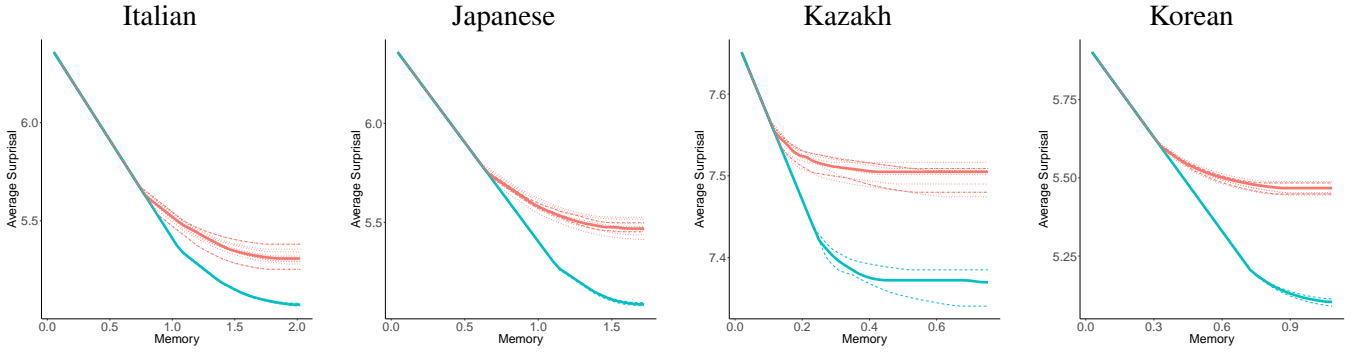
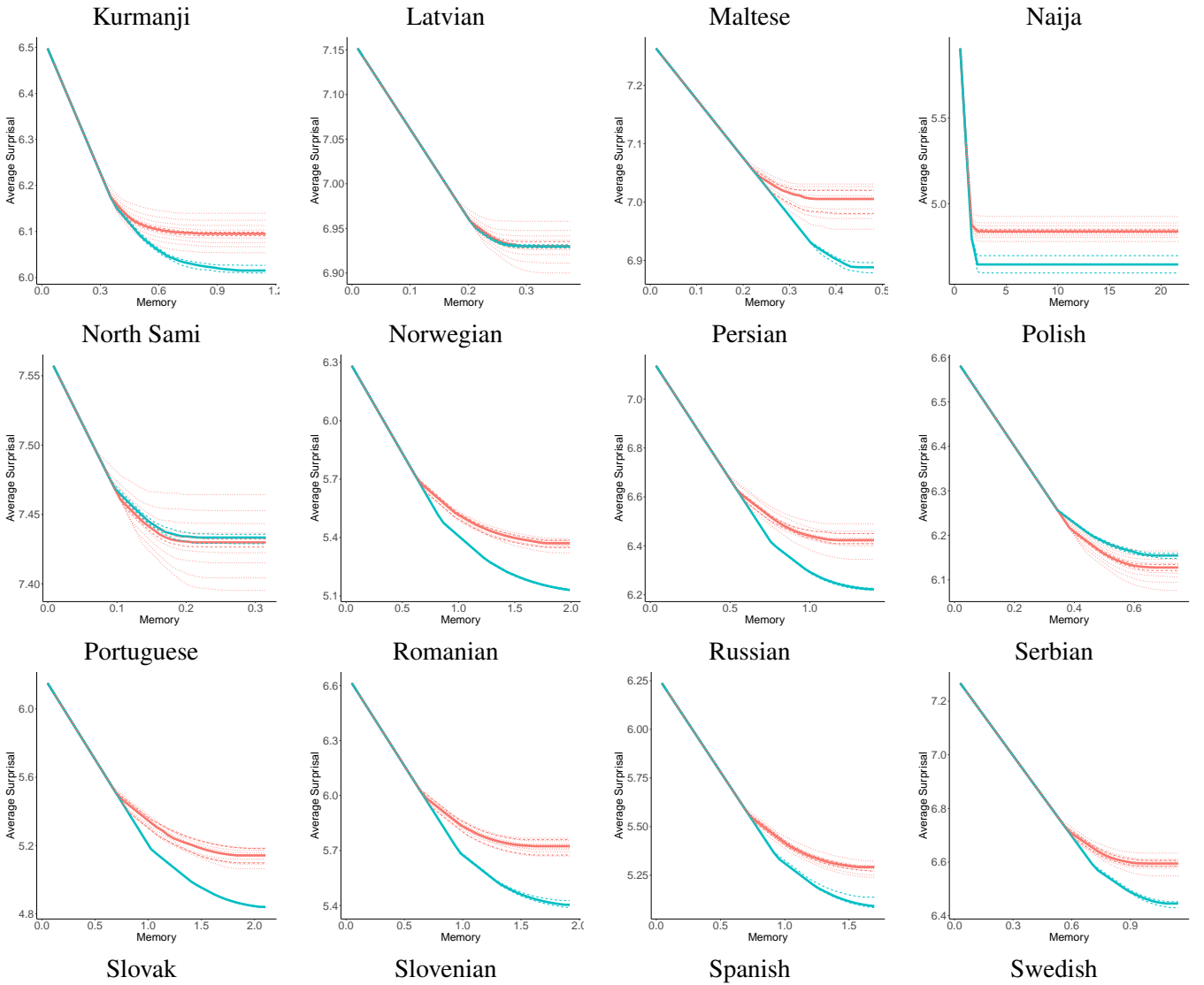


Figure 4: Medians (cont.)



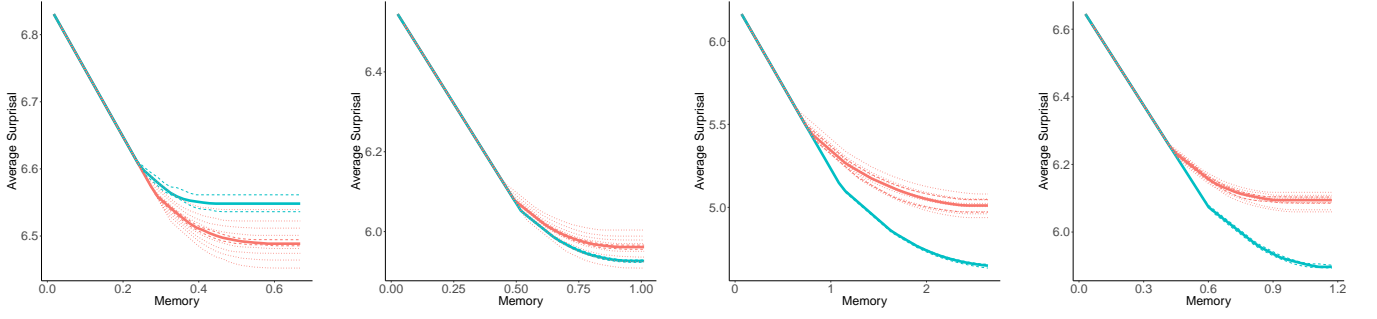


Figure 5: Medians (cont.)

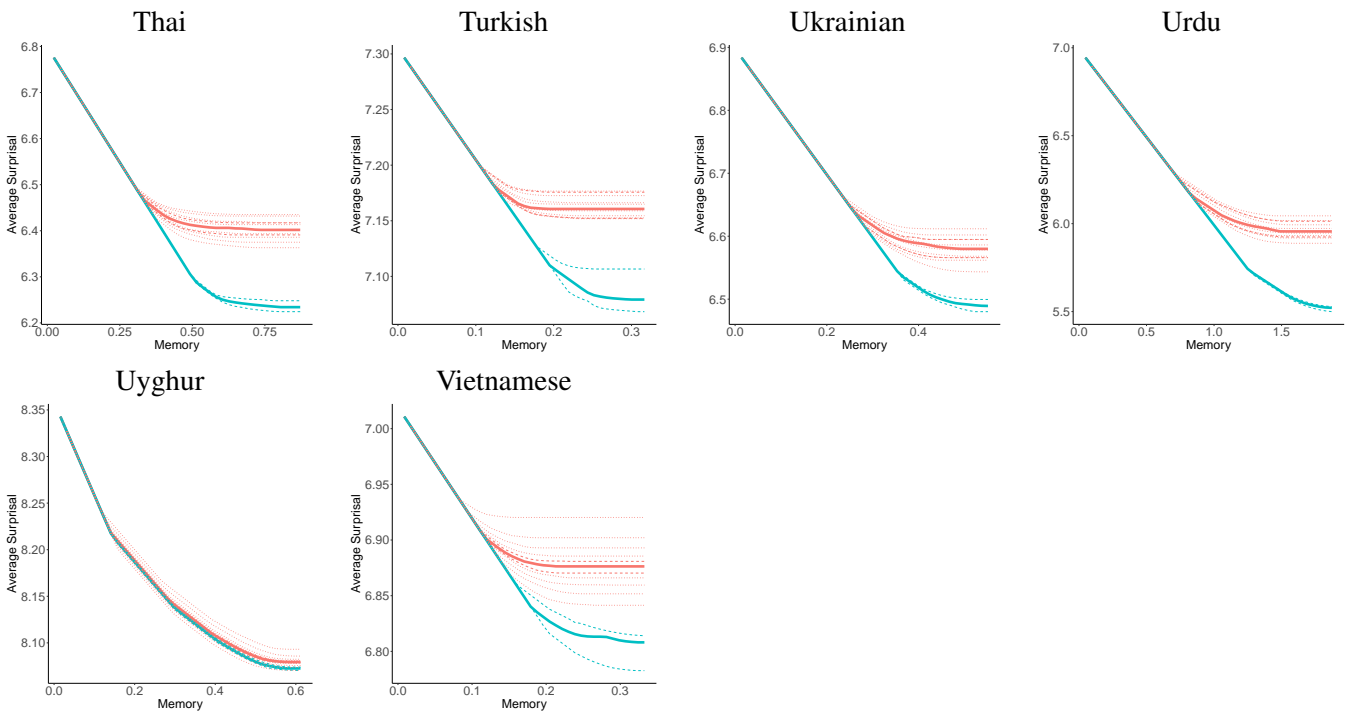


Figure 6: Medians (cont.)

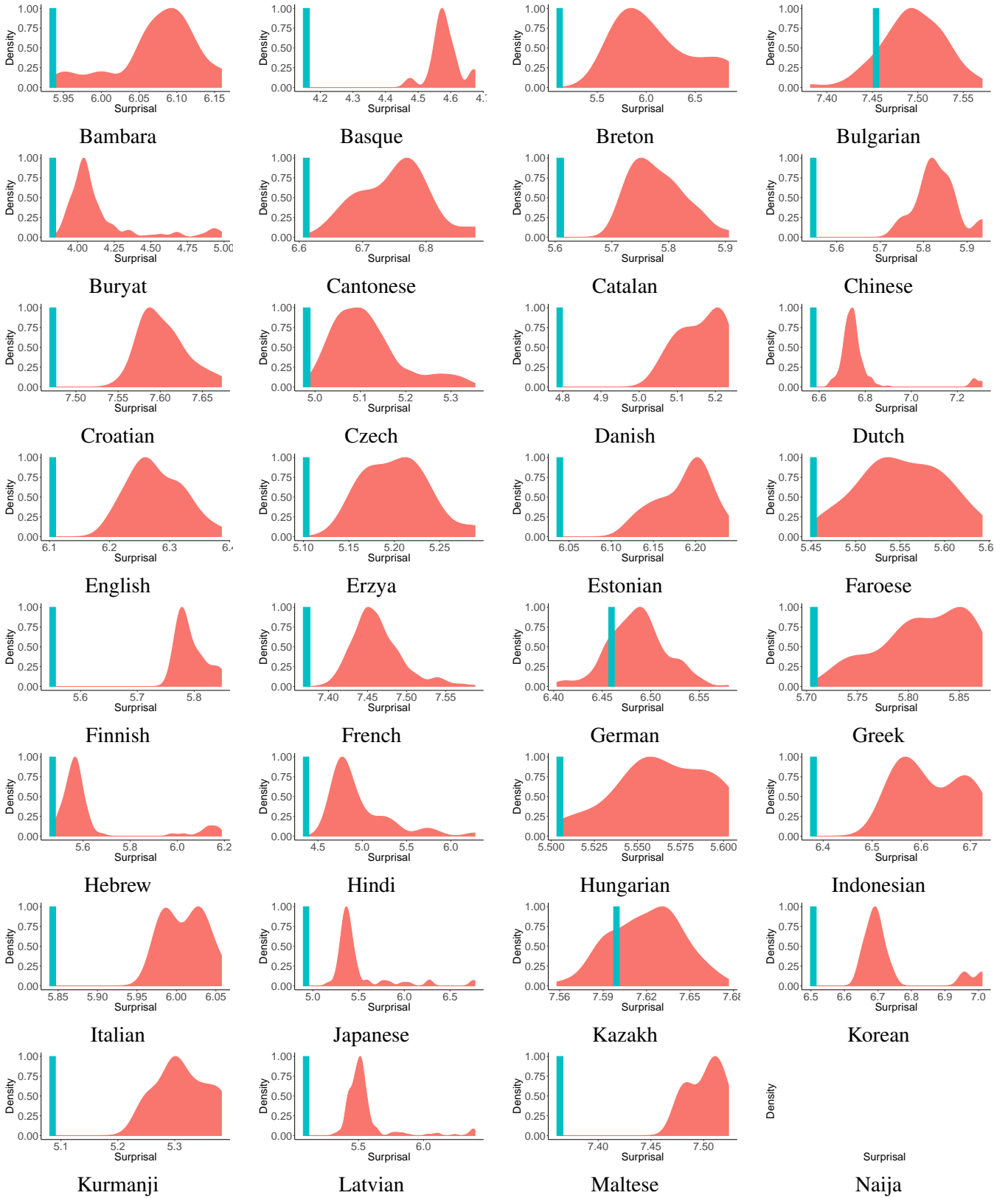
## 6.2 Surprisal at Maximum Memory

Afrikaans

Amharic

Arabic

Armenian





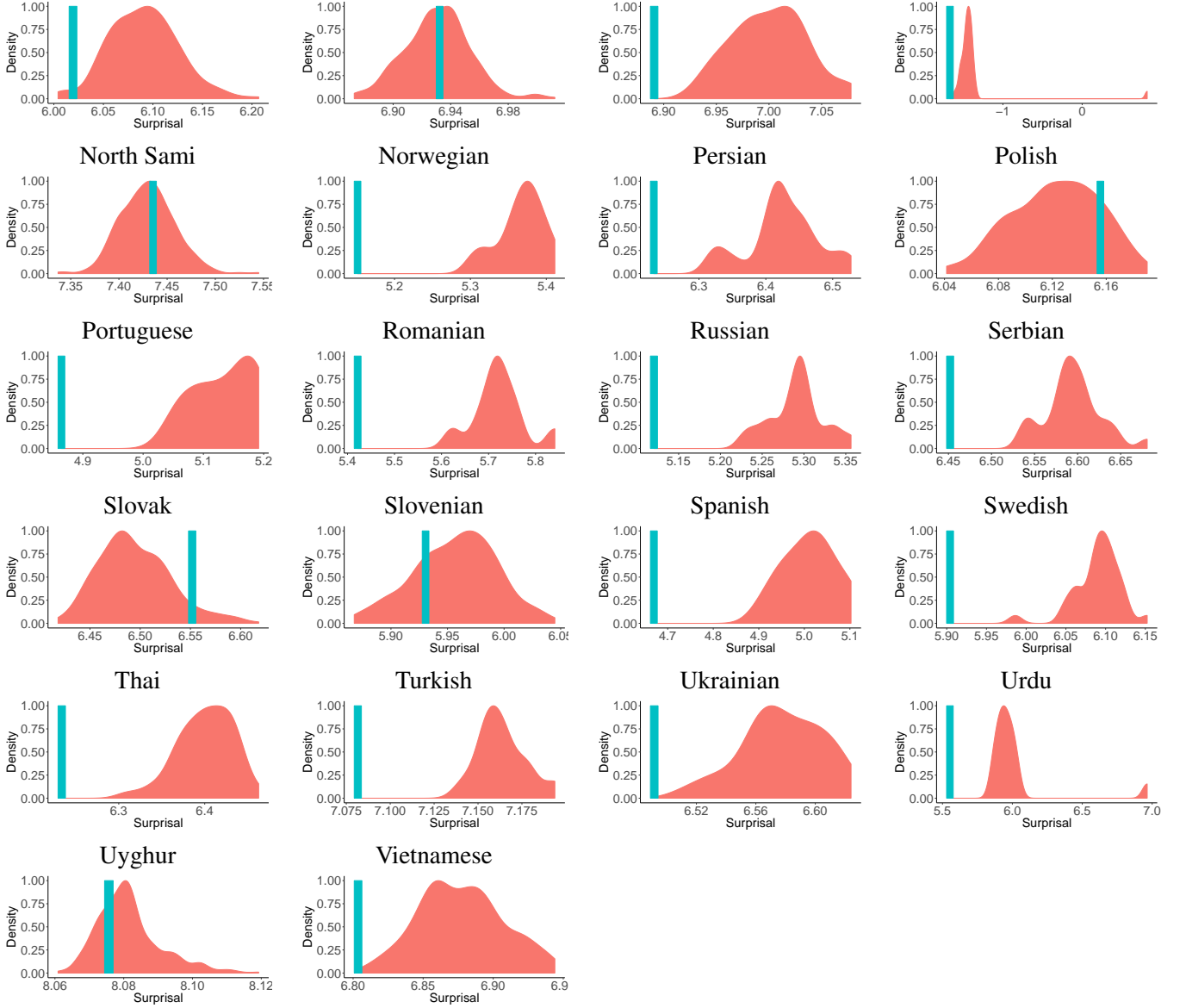


Figure 7: Histograms: Surprisal, at maximum memory.

{tab:slice-h

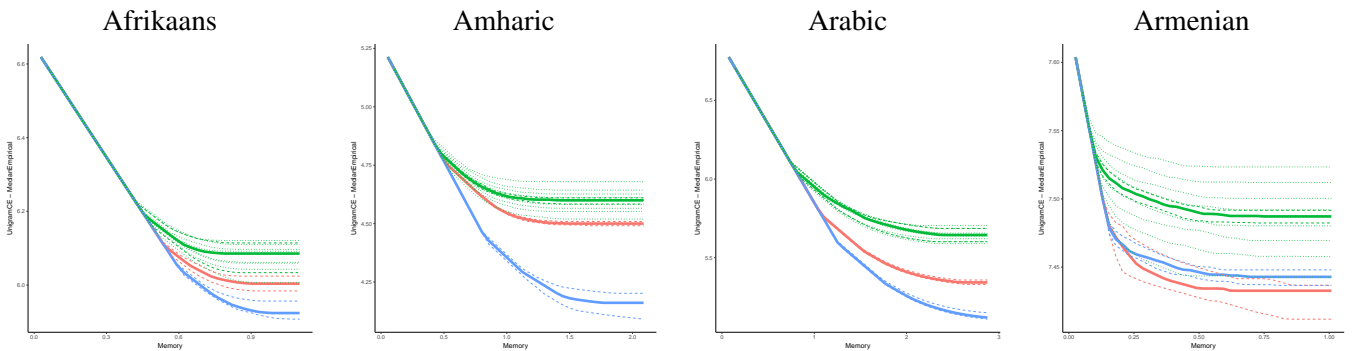
### 6.3 Samples Drawn (Experiment 3)

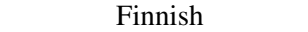
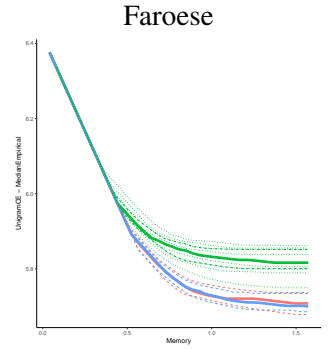
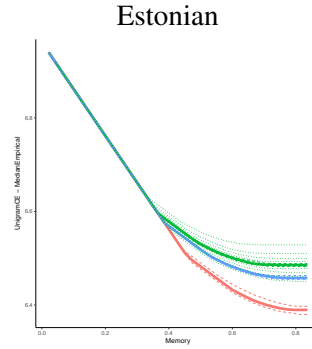
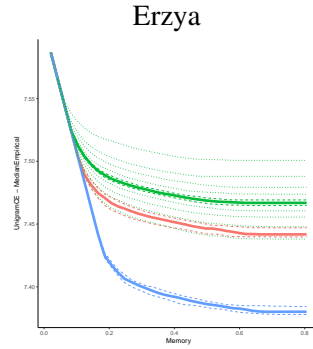
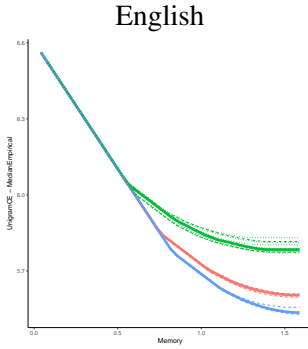
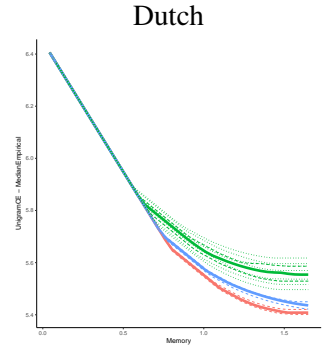
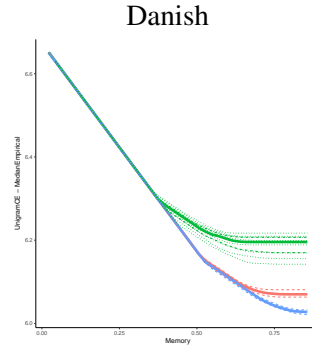
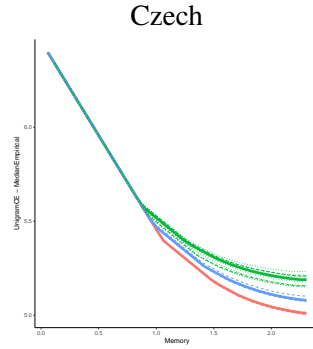
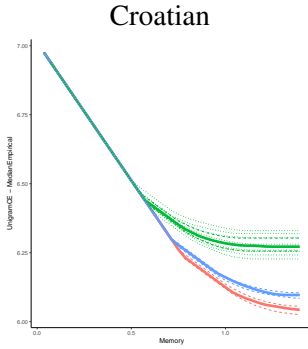
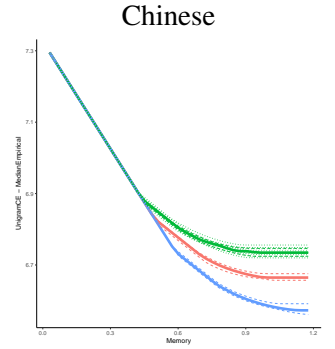
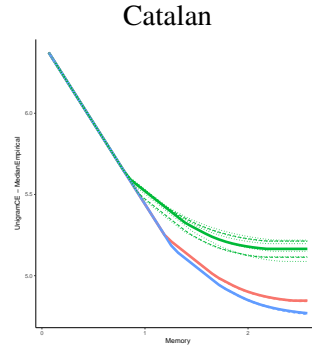
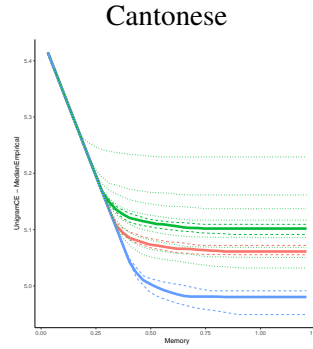
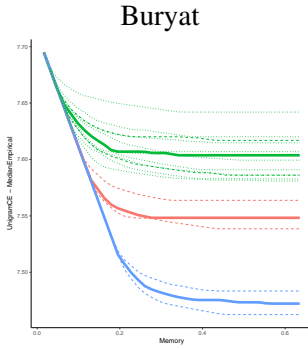
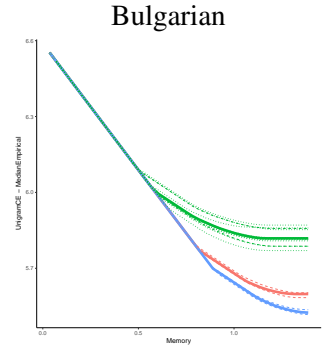
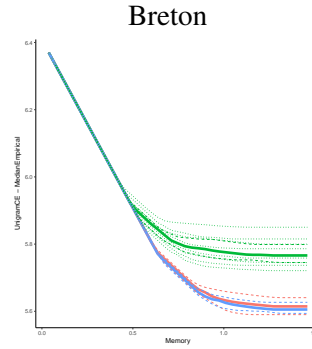
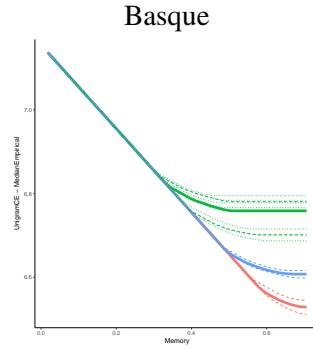
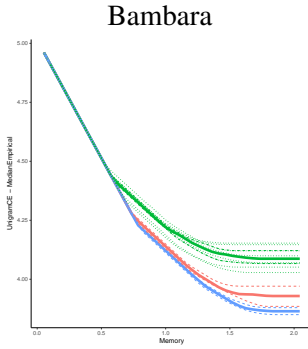
Language	Base.	MLE	Language	Base.	MLE
Afrikaans	13	10	Indonesian	11	10
Amharic	137	71	Italian	10	10
Arabic	11	10	Japanese	25	10
Armenian	140	17	Kazakh	11	10
Bambara	25	10	Korean	11	10
Basque	15	10	Kurmanji	338	101
Breton	35	10	Latvian	308	132
Bulgarian	14	10	Maltese	30	10
Buryat	26	10	Naija	214	93
Cantonese	306	135	North Sami	335	101
Catalan	11	10	Norwegian	12	10
Chinese	21	10	Persian	25	10
Croatian	30	10	Polish	309	131
Czech	18	12	Portuguese	15	99
Danish	33	10	Romanian	10	10
Dutch	27	10	Russian	20	13
English	13	10	Serbian	26	11
Erzya	846	101	Slovak	303	138
Estonian	347	10	Slovenian	297	12
Faroese	27	10	Spanish	14	10
Finnish	83	54	Swedish	31	10
French	14	12	Thai	45	10
German	19	10	Turkish	13	10
Greek	16	10	Ukrainian	28	10
Hebrew	11	10	Urdu	17	10
Hindi	11	10	Uyghur	326	132
Hungarian	220	35	Vietnamese	303	132

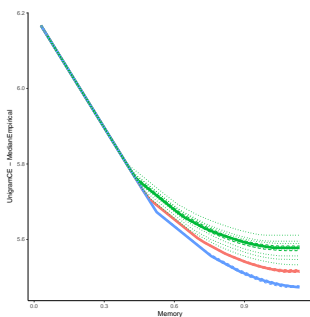
Figure 8: Experiment 3: Samples drawn per language according to the precision-dependent stopping criterion.

{tab:samples}

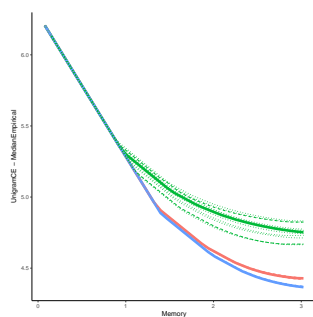
### 6.4 Medians (Experiment 3)



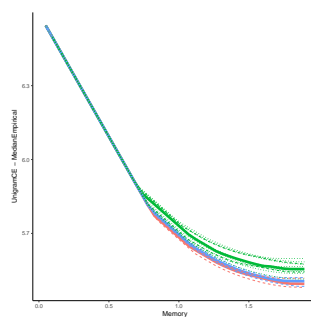




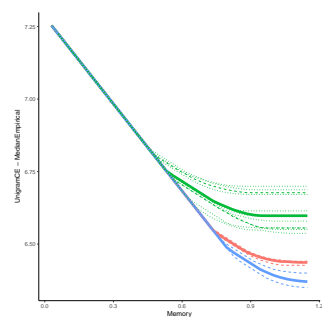
Hebrew



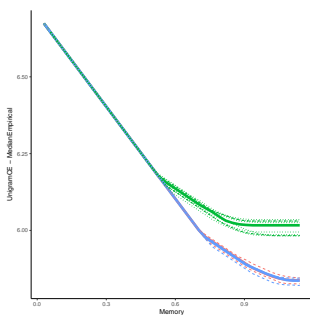
Hindi



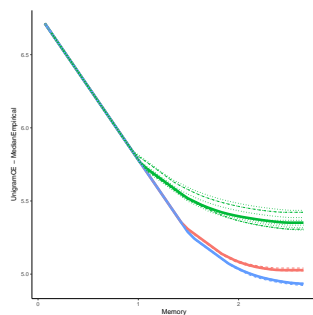
Hungarian



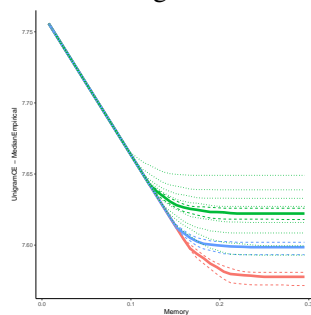
Indonesian



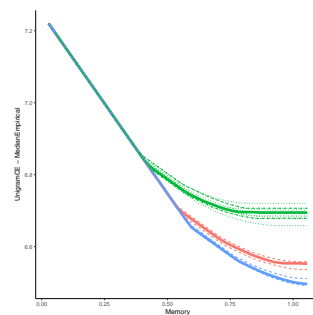
Italian



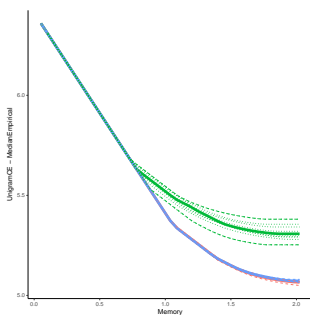
Japanese



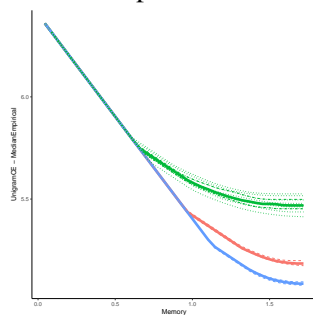
Kazakh



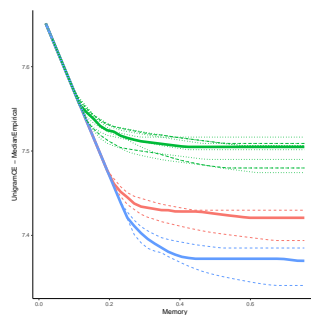
Korean



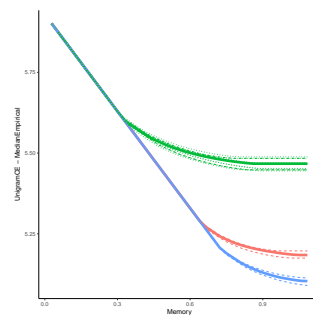
Kurmanji



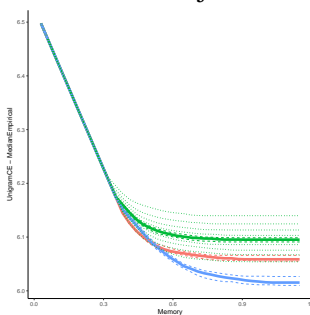
Latvian



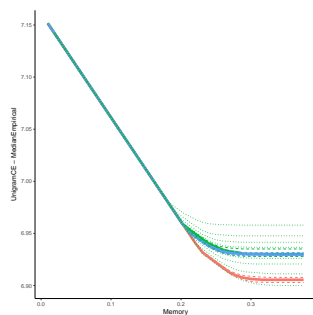
Maltese



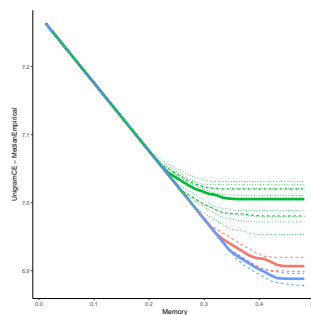
Naija



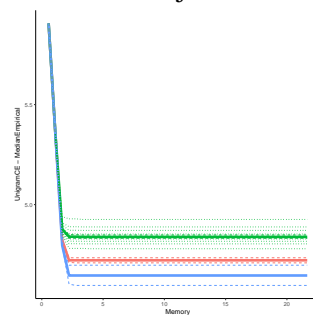
North Sami



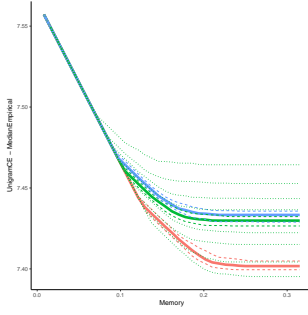
Norwegian



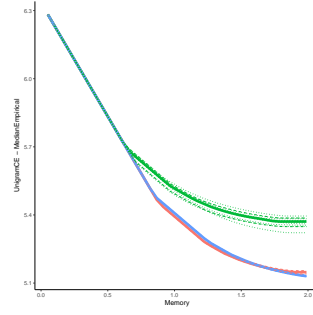
Persian



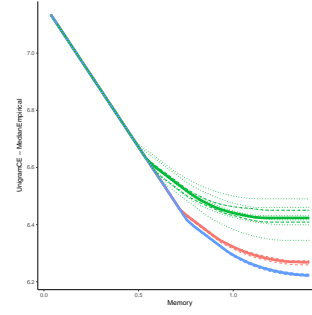
Polish



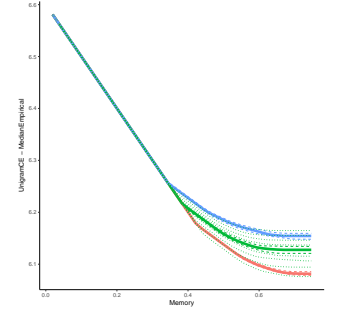
Portuguese



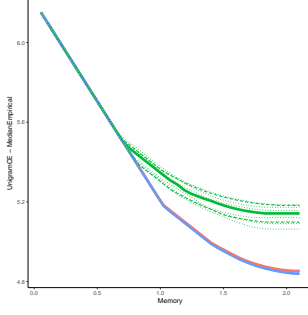
Romanian



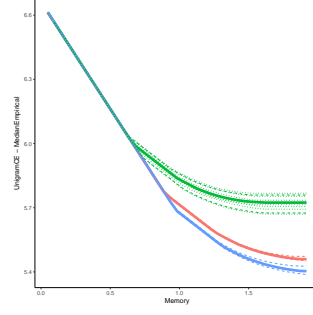
Russian



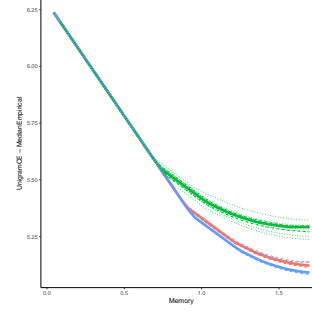
Serbian



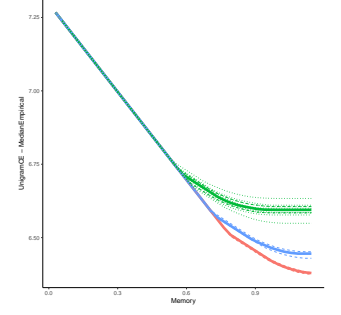
Slovak



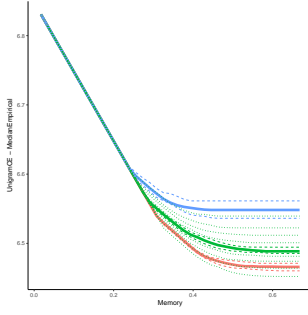
Slovenian



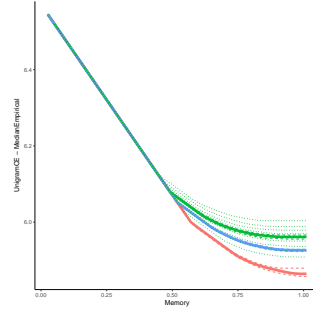
Spanish



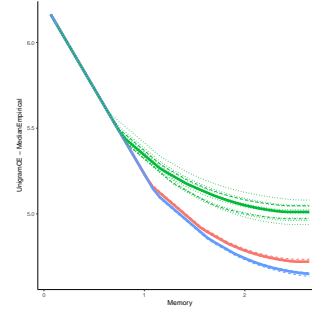
Swedish



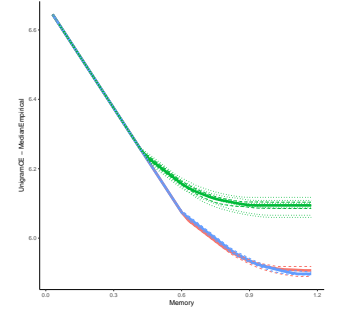
Thai



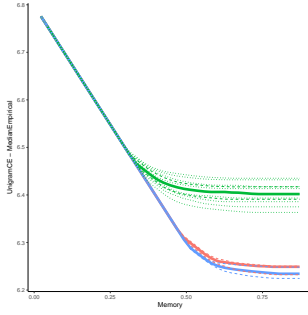
Turkish



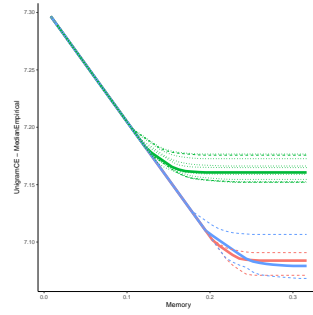
Ukrainian



Urdu



Uyghur



Vietnamese

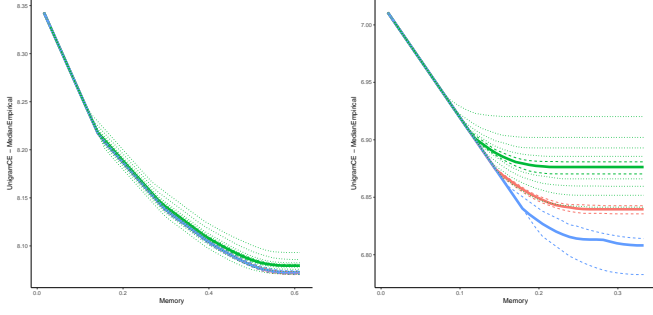
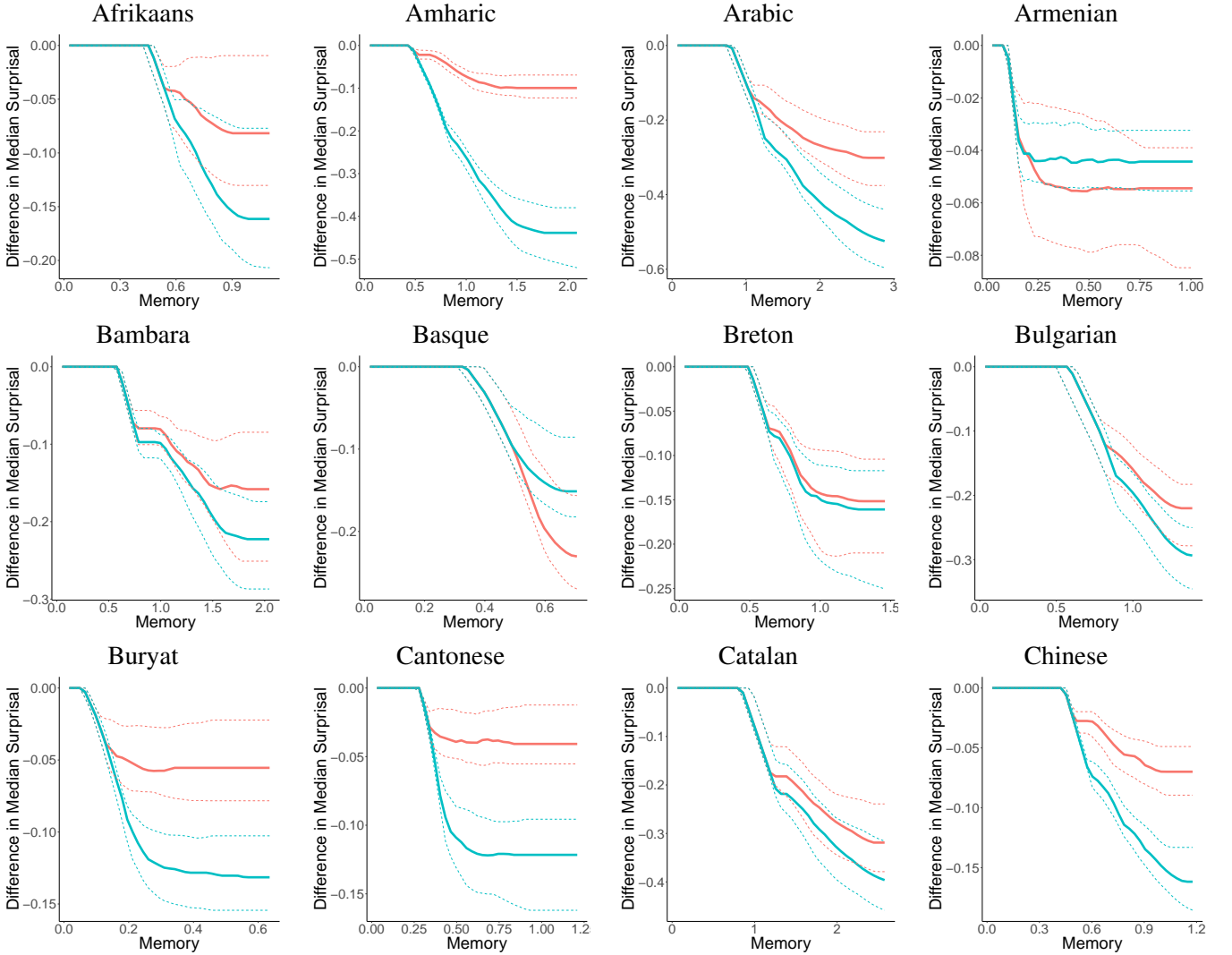
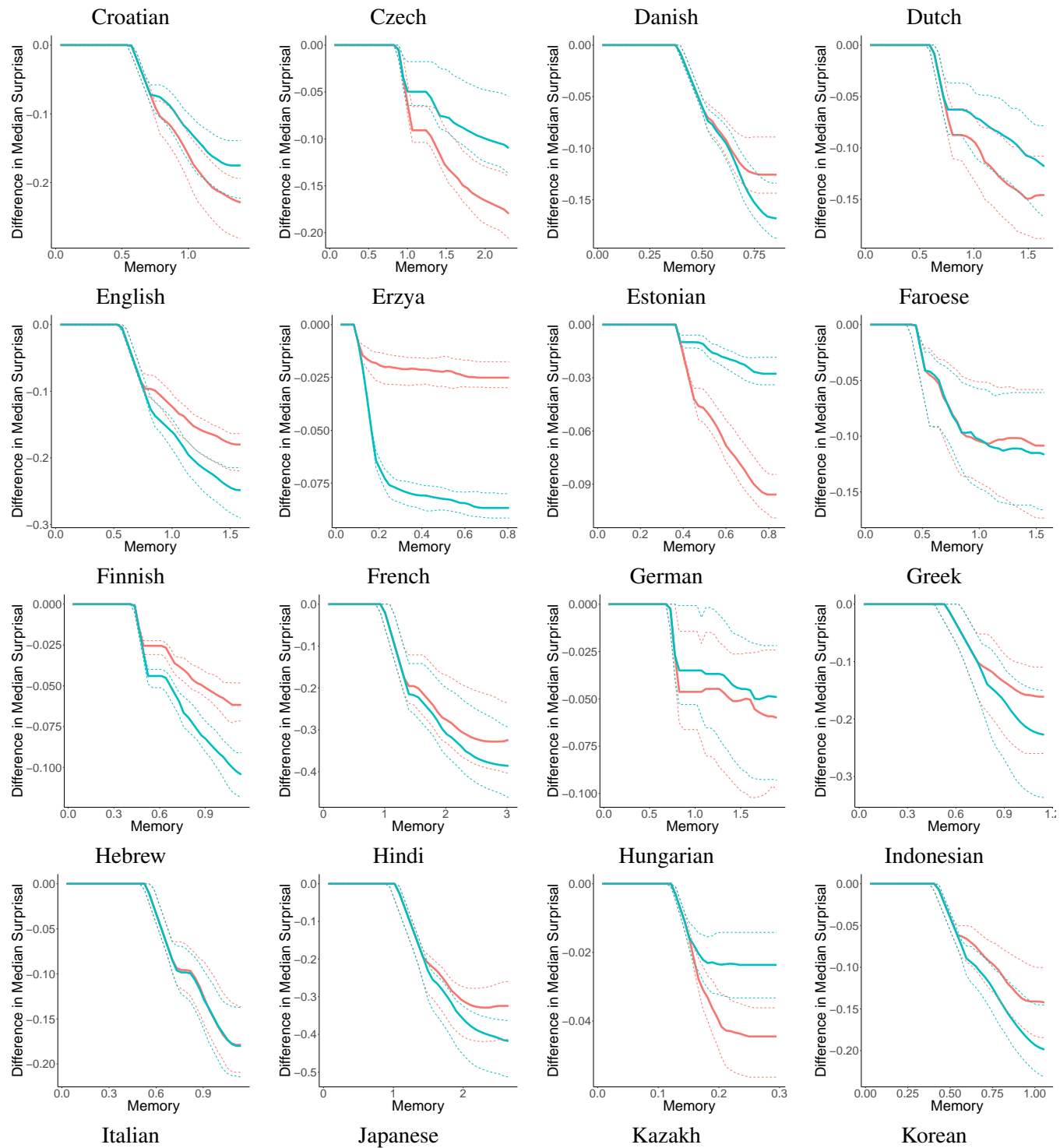
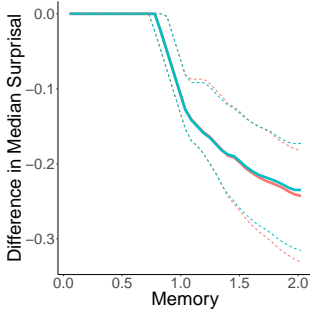


Figure 9: Experiment 3. Medians: For each memory budget, we provide the median surprisal for real and random languages. Solid lines indicate sample medians, dashed lines indicate 95 % confidence intervals for the population median. Green: Random baselines; blue: real language; red: maximum-likelihood grammars fit to real orderings.

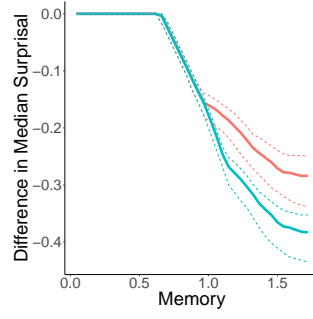
{tab:medians



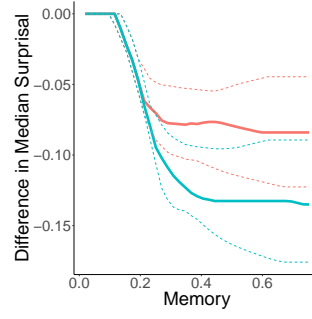




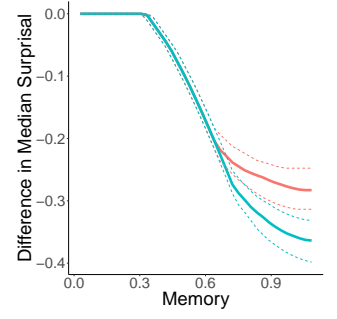
Kurmanji



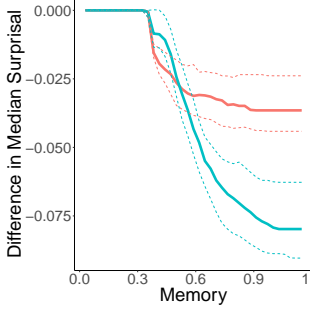
Latvian



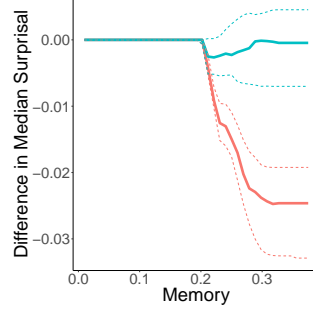
Maltese



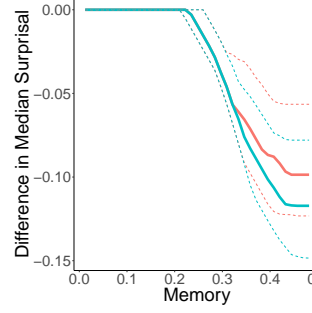
Naija



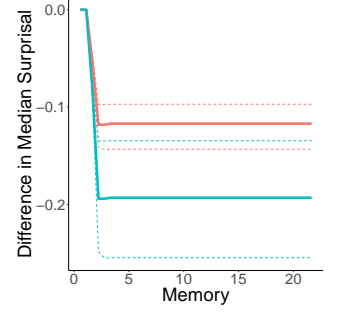
North Sami



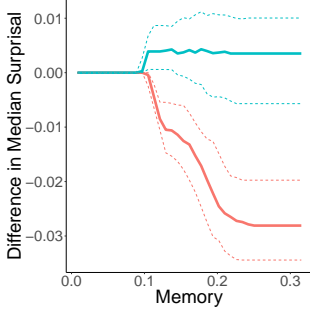
Norwegian



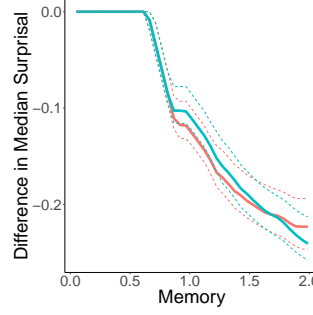
Persian



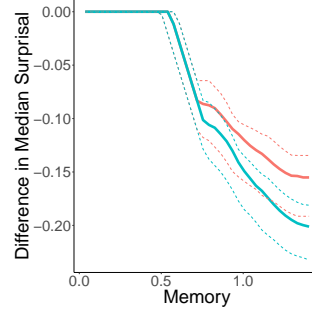
Polish



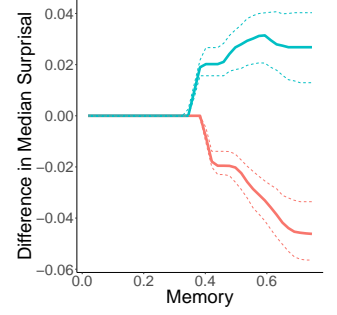
Portuguese



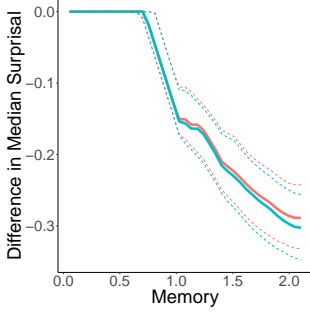
Romanian



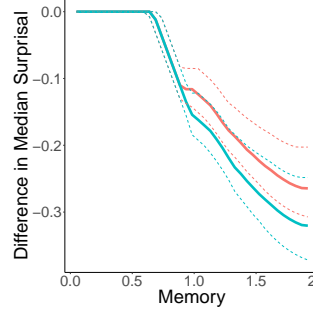
Russian



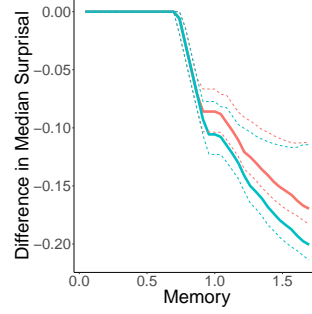
Serbian



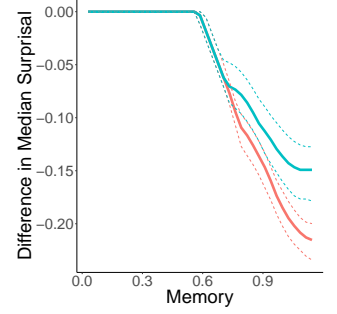
Slovak



Slovenian



Spanish



Swedish



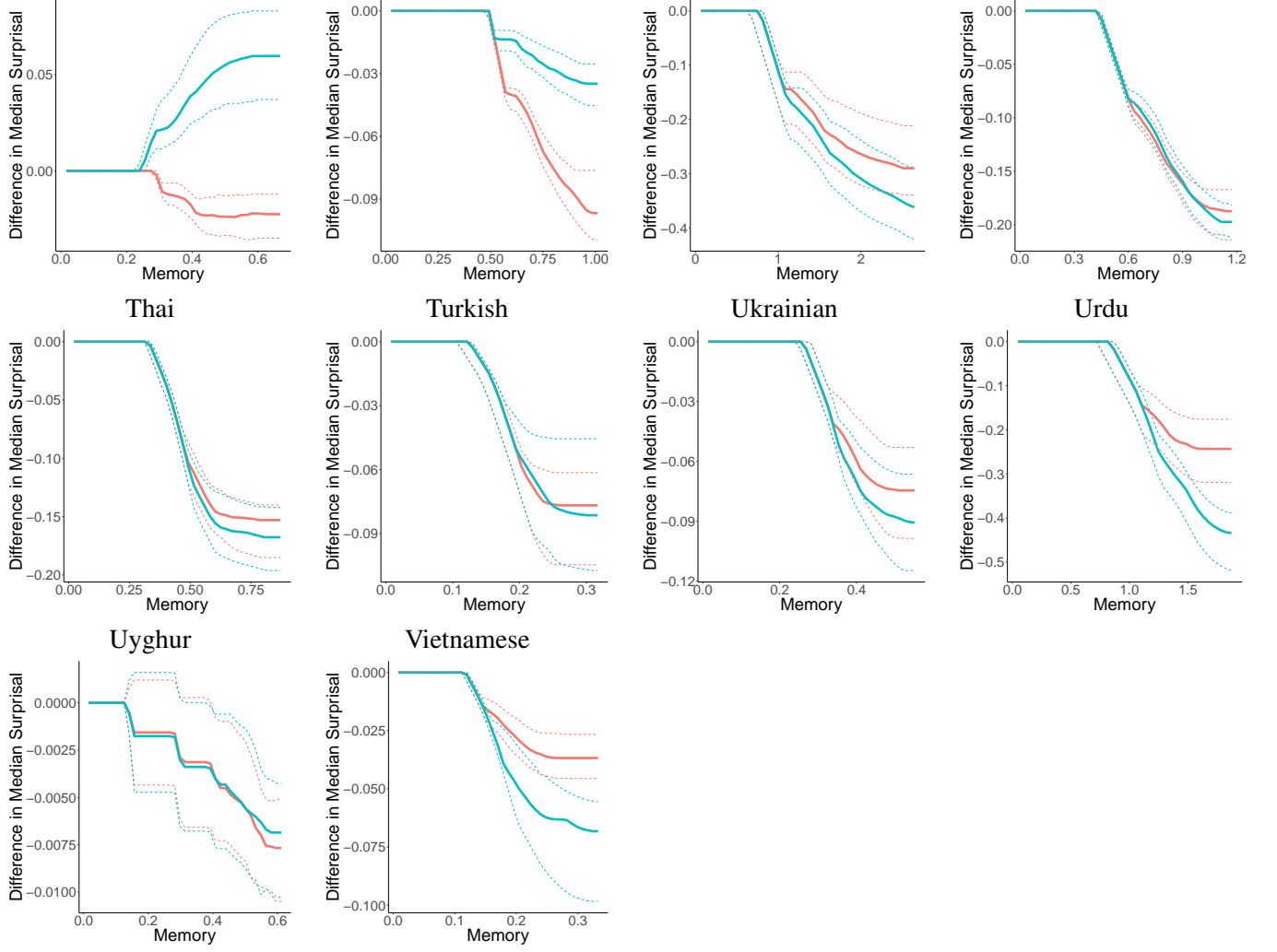
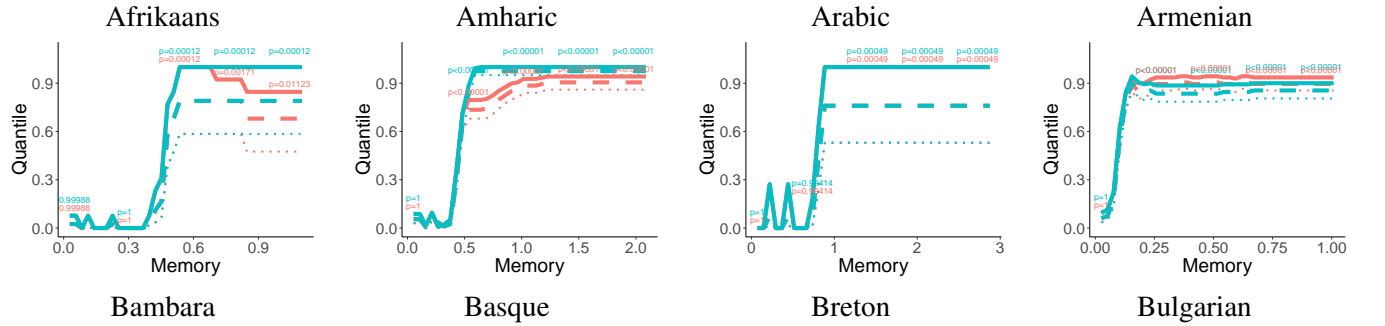
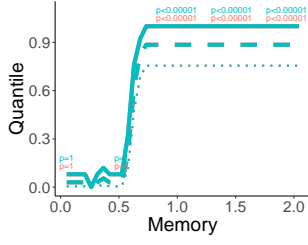


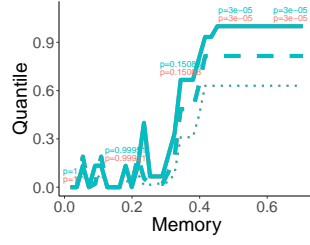
Figure 10: Median Differences between Real and Baseline: For each memory budget, we provide the difference in median surprisal between real languages and random baselines; for real orders (blue) and maximum likelihood grammars (red). Lower values indicate lower surprisal compared to baselines. Solid lines indicate sample means. Dashed lines indicate 95 % confidence intervals.

{tab:median\_

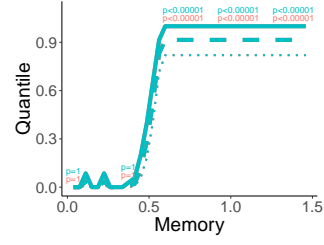




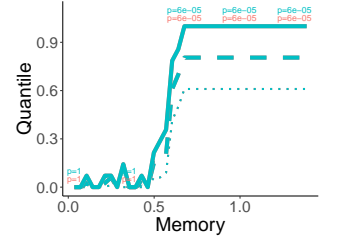
Buryat



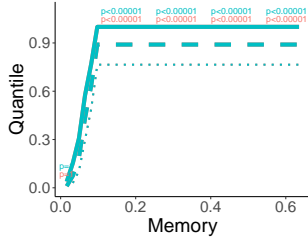
Cantonese



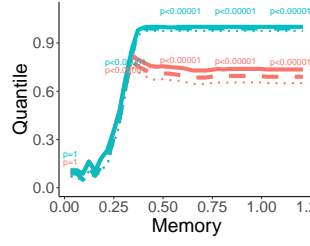
Catalan



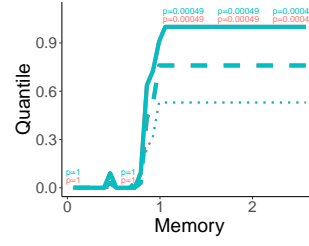
Chinese



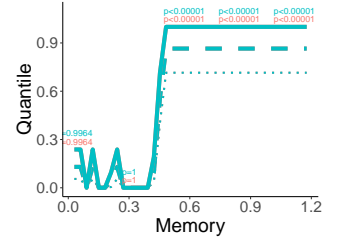
Croatian



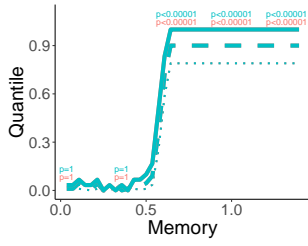
Czech



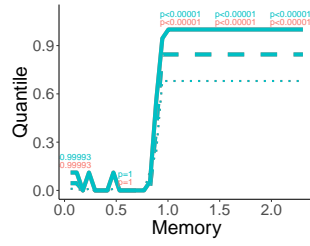
Danish



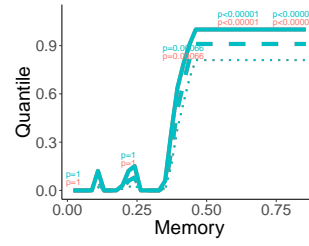
Dutch



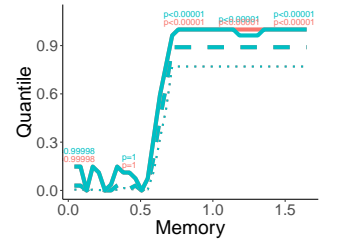
English



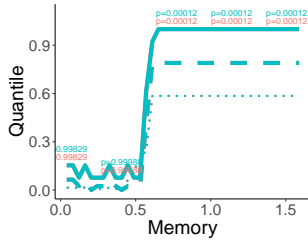
Erzya



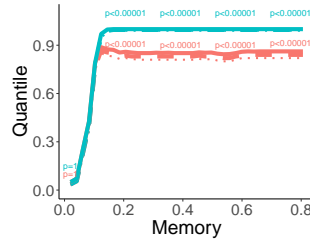
Estonian



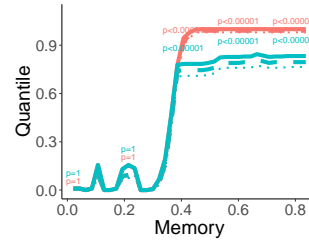
Faroese



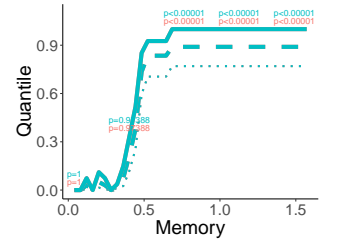
Finnish



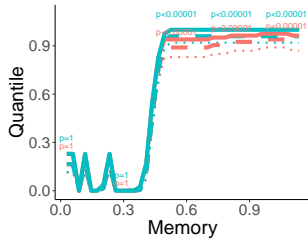
French



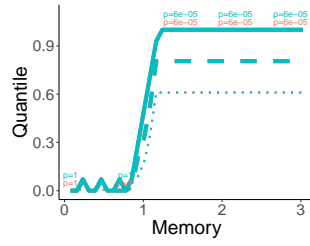
German



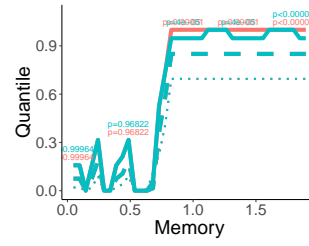
Greek



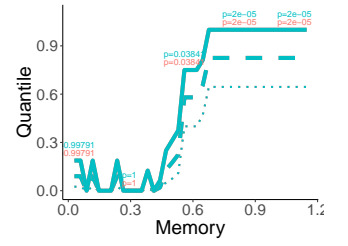
Hebrew



Hindi



Hungarian



Indonesian



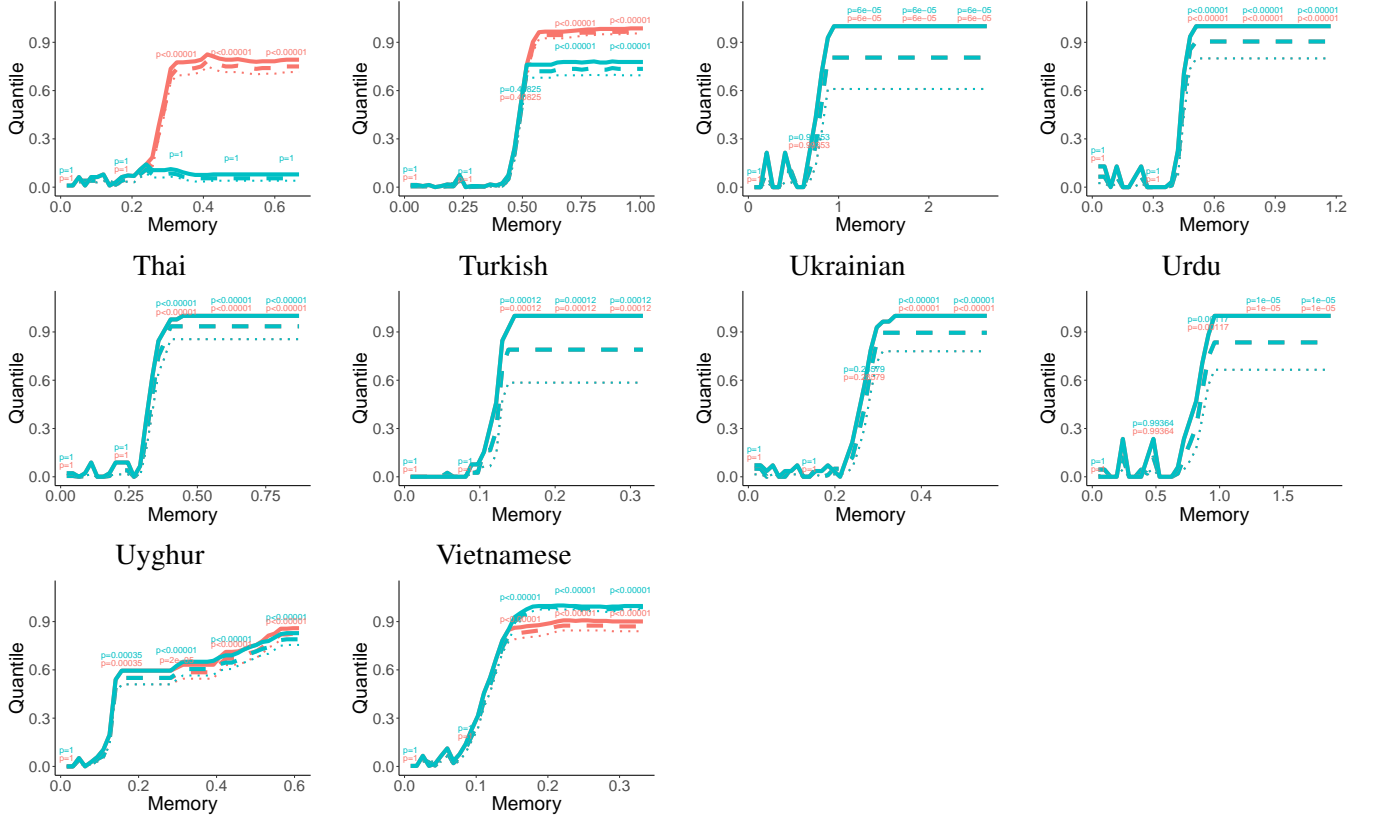


Figure 11: Quantiles: At a given memory budget, what percentage of the baselines results in higher listener surprisal than the real language? Solid curves represent sample means, dashed lines represent 95 % confidence bounds; dotted lines represent 99.9 % confidence bounds. At five evenly spaced memory levels, we provide a p-value for the null hypothesis that the actual population mean is 0.5 or less. Confidence bounds and p-values are obtained using an exact nonparametric method (see text).

{tab:quantil

## 7 Details for Neural Network Models

## 8 N-Gram Models

### 8.1 Method

We use a version of Kneser-Ney Smoothing. For a sequence  $w_1 \dots w_k$ , let  $N(w_{1\dots k})$  be the number of times  $w_{1\dots k}$  occurs in the training set. The unigram probabilities are estimated as

$$p_1(w_t) := \frac{N(w_t) + \delta}{|Train| + |V| \cdot \delta} \quad (47)$$

where  $\delta \in \mathbb{R}_+$  is a hyperparameter. Here  $|Train|$  is the number of tokens in the training set,  $|V|$  is the number of types occurring in train or held-out data. Higher-order probabilities  $p_t(w_t|w_{0\dots t-1})$  are estimated

recursively as follows. Let  $\gamma > 0$  be a hyperparameter. If  $N(w_{0\dots t-1}) < \gamma$ , set

$$p_t(w_t|w_{0\dots t-1}) := p_{t-1}(w_t|w_{1\dots t-1}) \quad (48)$$

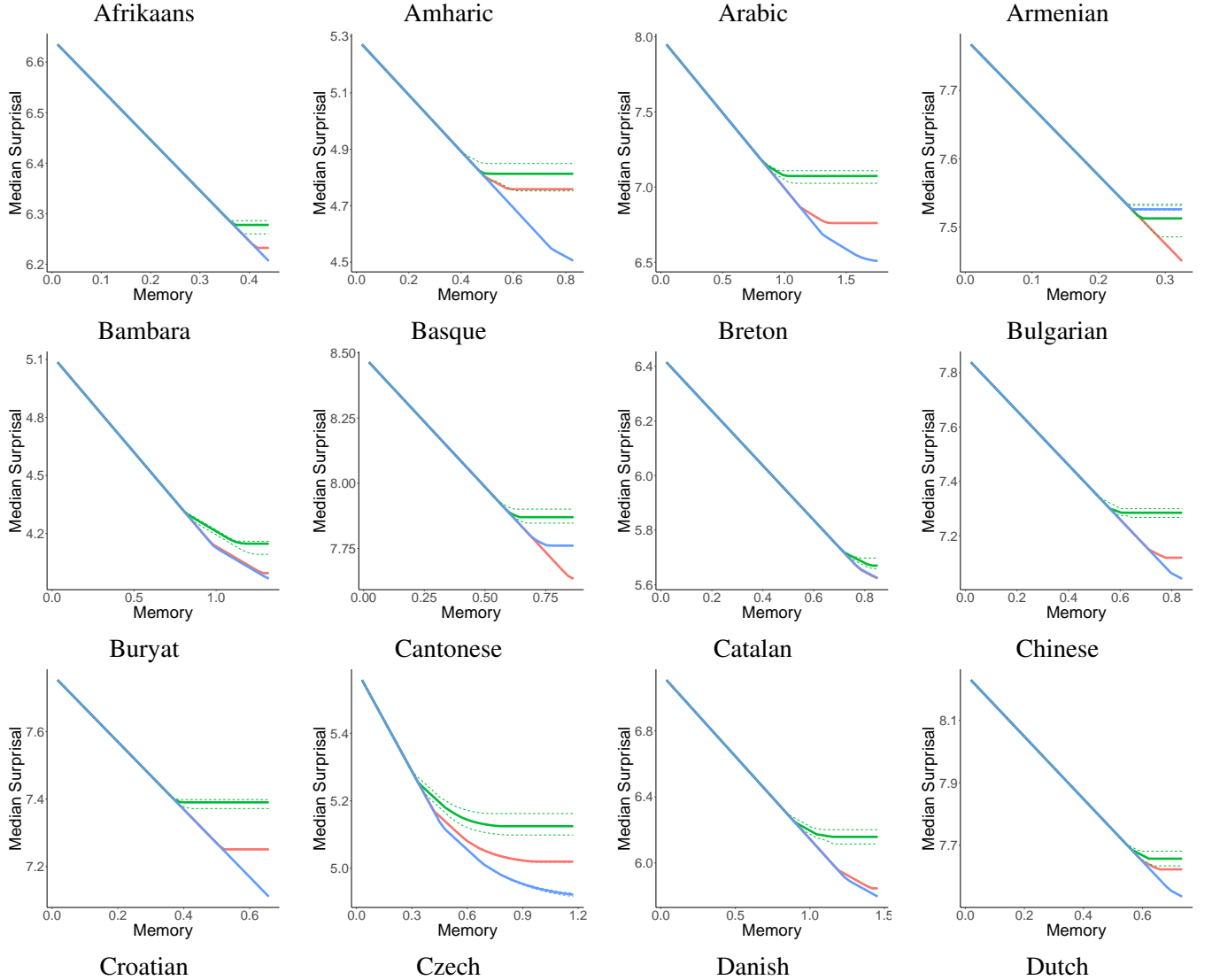
Otherwise, we interpolate between  $t$ -th order and lower-order estimates:

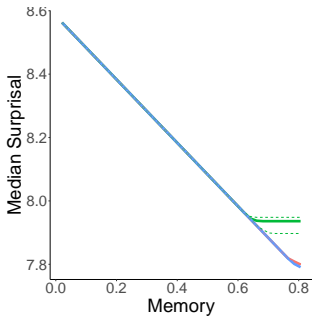
$$p_t(w_t|w_{0\dots t-1}) := \frac{\max(N(w_{0\dots t}) - \alpha, 0.0) + \alpha \cdot \#\{w : N(w_{0\dots t-1}w) > 0\} \cdot p_{t-1}(w_t|w_{1\dots t-1})}{N(w_{0\dots t-1})} \quad (49)$$

where  $\alpha \in [0, 1]$  is also a hyperparameter. (CITE) show that this definition results in a well-defined probability distribution, i.e.,  $\sum_{w \in V} p_t(w|w_{0\dots t-1}) = 1$ .

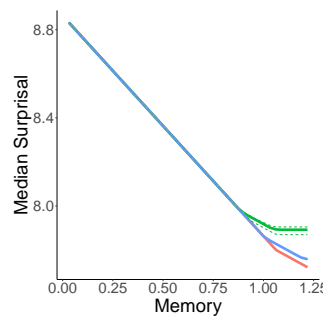
Hyperparameters  $\alpha, \gamma, \delta$  are tuned with the same strategy as for the neural network models.

## 8.2 Results

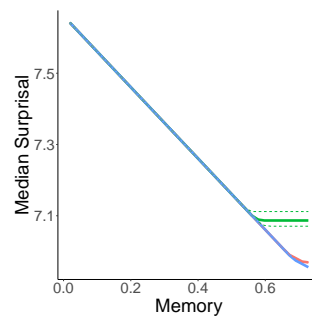




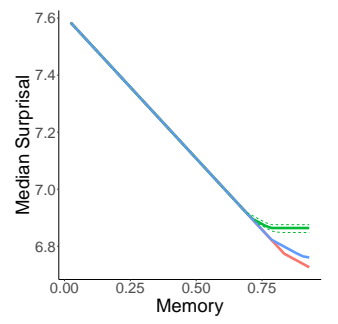
English



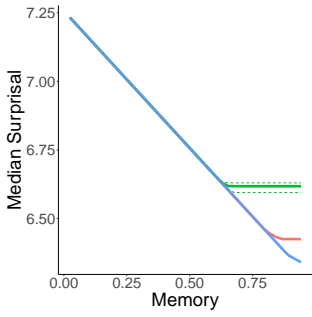
Erzya



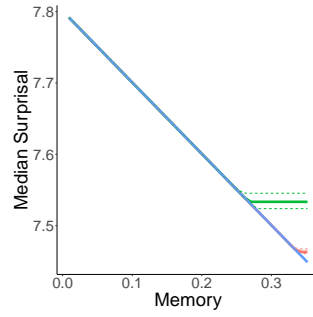
Estonian



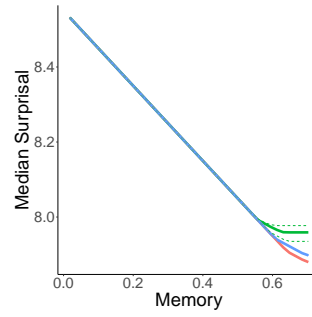
Faroese



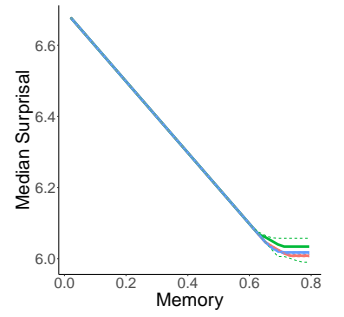
Finnish



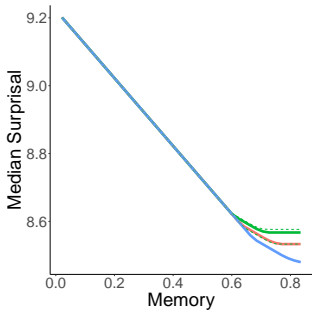
French



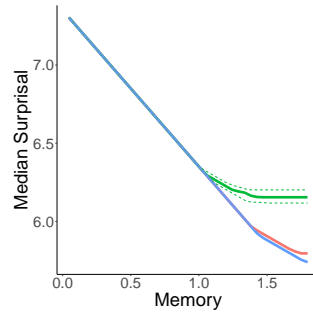
German



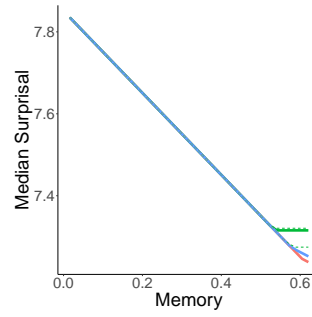
Greek



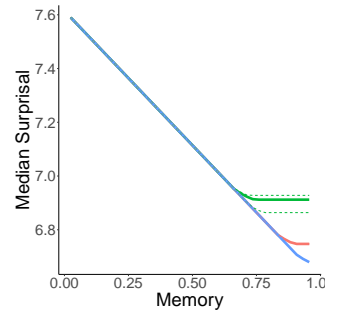
Hebrew



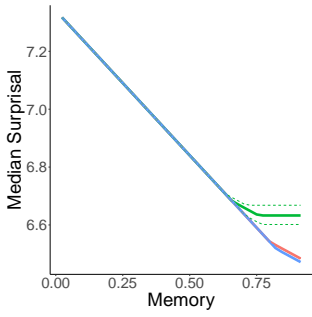
Hindi



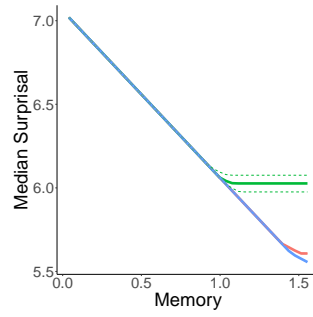
Hungarian



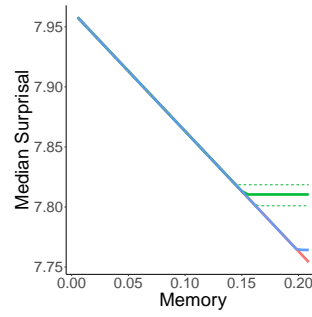
Indonesian



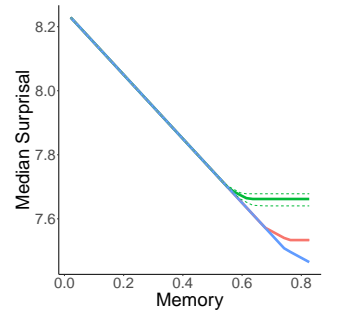
Italian



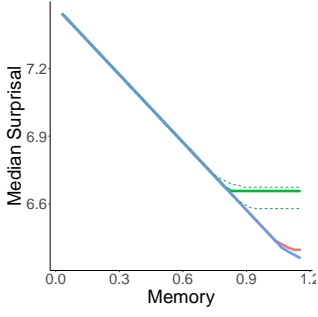
Japanese



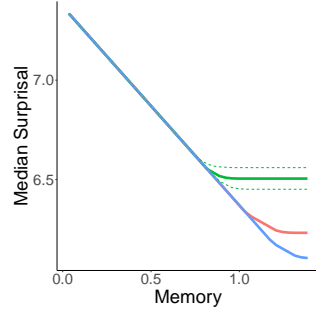
Kazakh



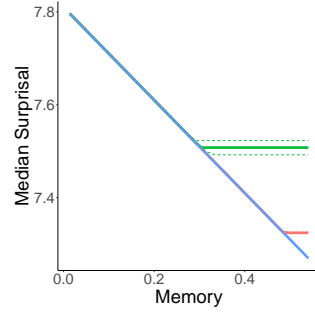
Korean



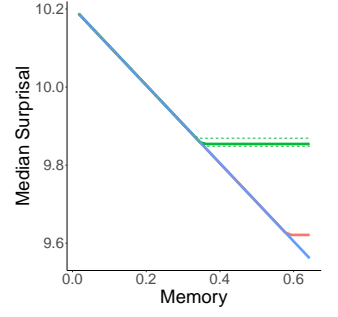
Kurmanji



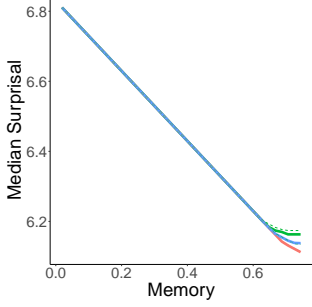
Latvian



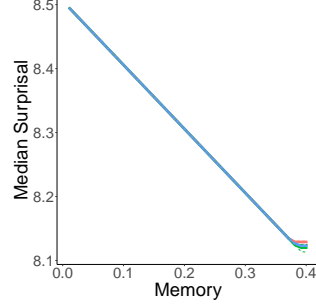
Maltese



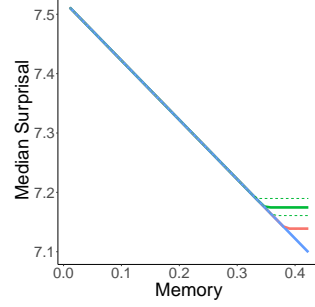
Naija



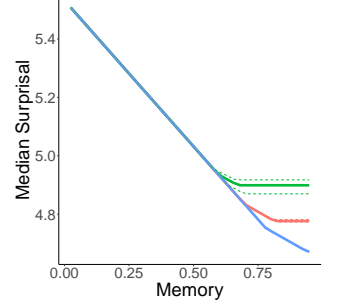
North Sami



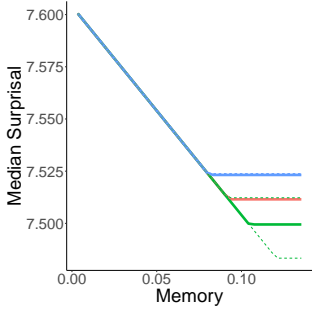
Norwegian



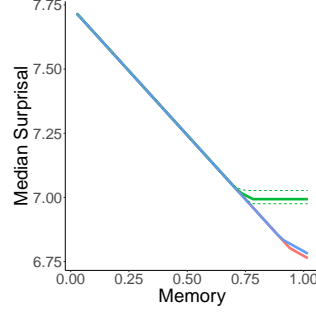
Persian



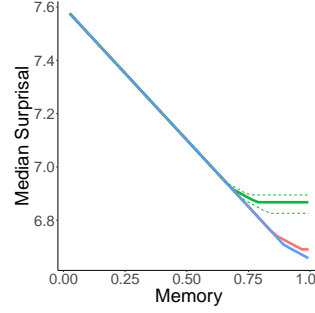
Polish



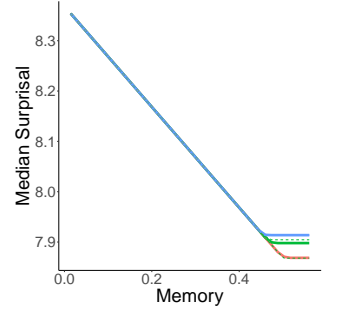
Portuguese



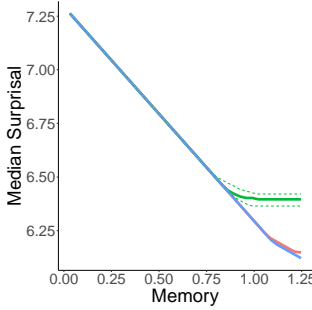
Romanian



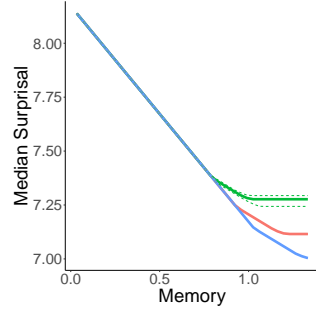
Russian



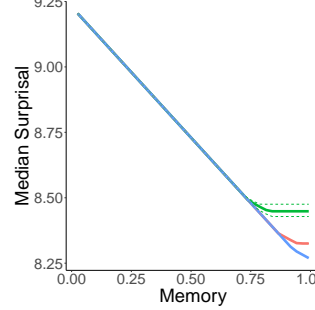
Serbian



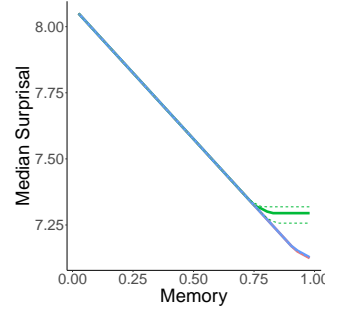
Slovak



Slovenian



Spanish



Swedish

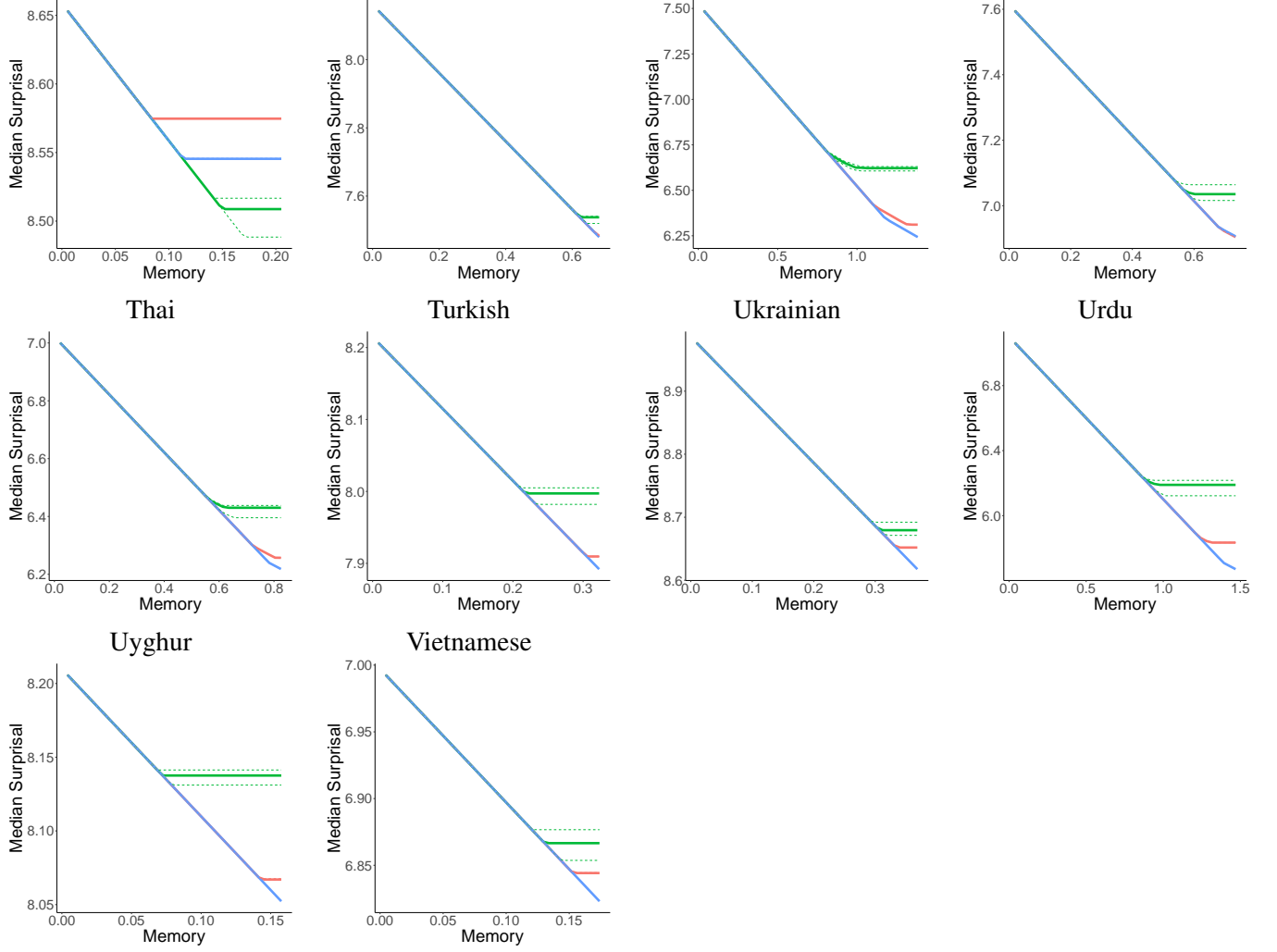
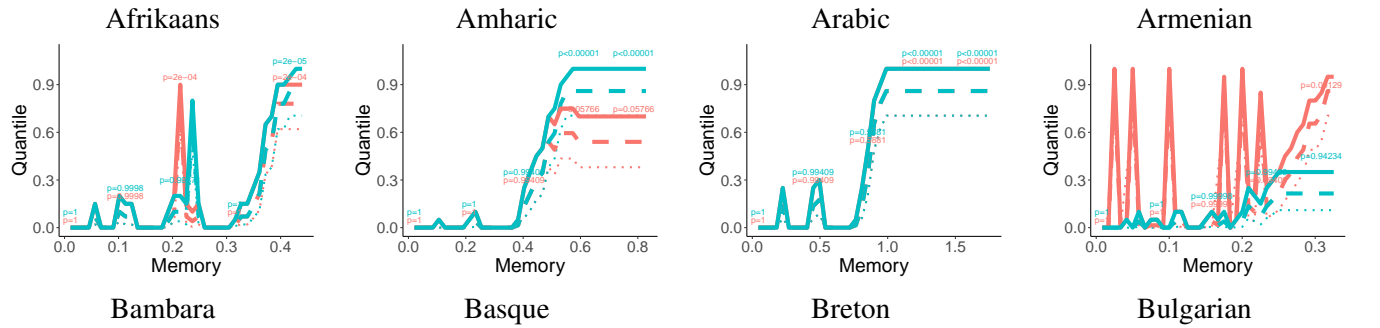
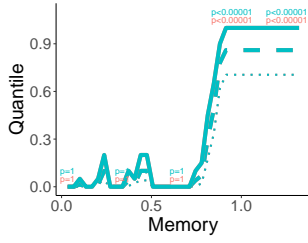


Figure 12: Medians (estimated using n-gram models): For each memory budget, we provide the median surprisal for real and random languages. Solid lines indicate sample medians for ngrams, dashed lines indicate 95 % confidence intervals for the population median. Green: Random baselines; blue: real language; red: maximum-likelihood grammars fit to real orderings.

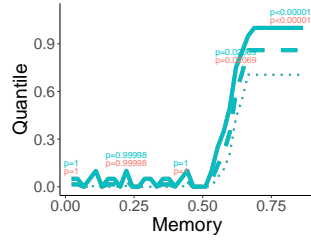
{tab:medians}



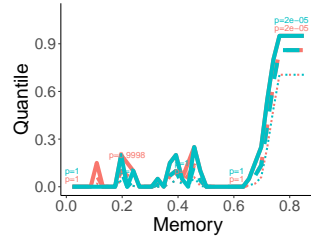




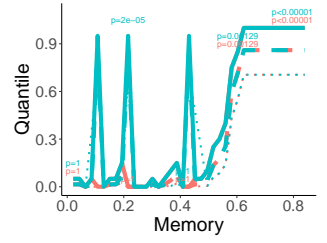
Buryat



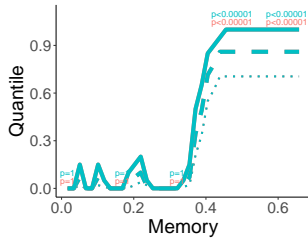
Cantonese



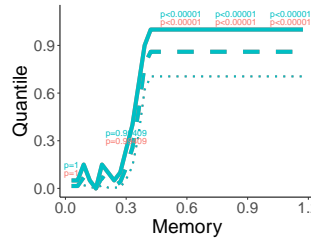
Catalan



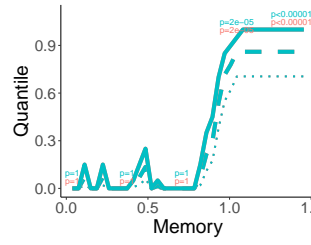
Chinese



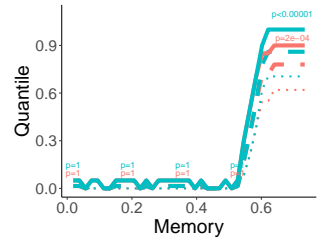
Croatian



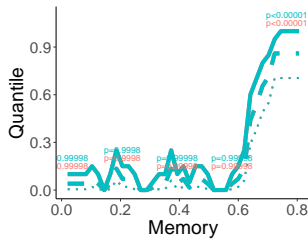
Czech



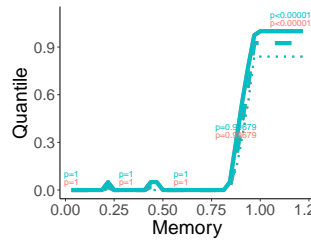
Danish



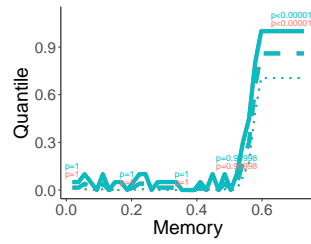
Dutch



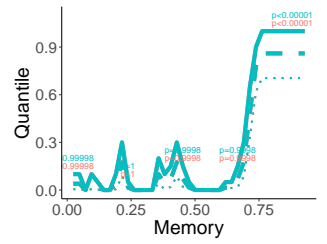
English



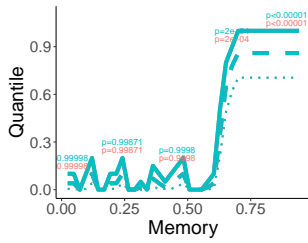
Erzya



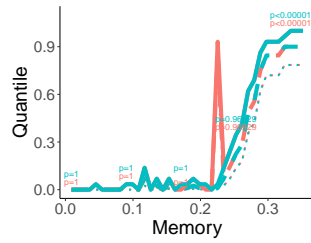
Estonian



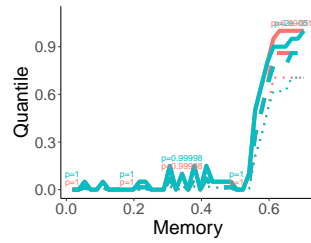
Faroese



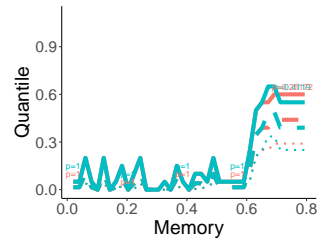
Finnish



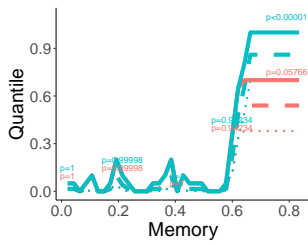
French



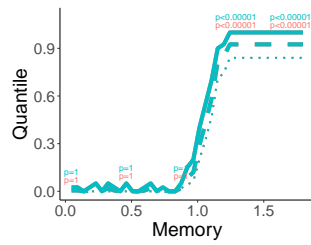
German



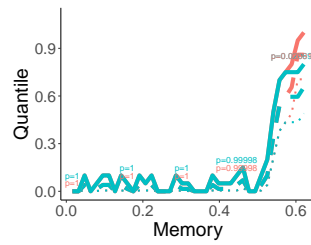
Greek



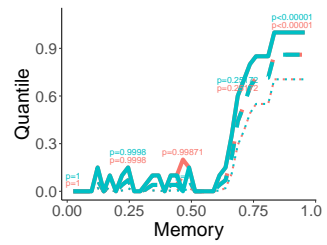
Hebrew



Hindi



Hungarian



Indonesian



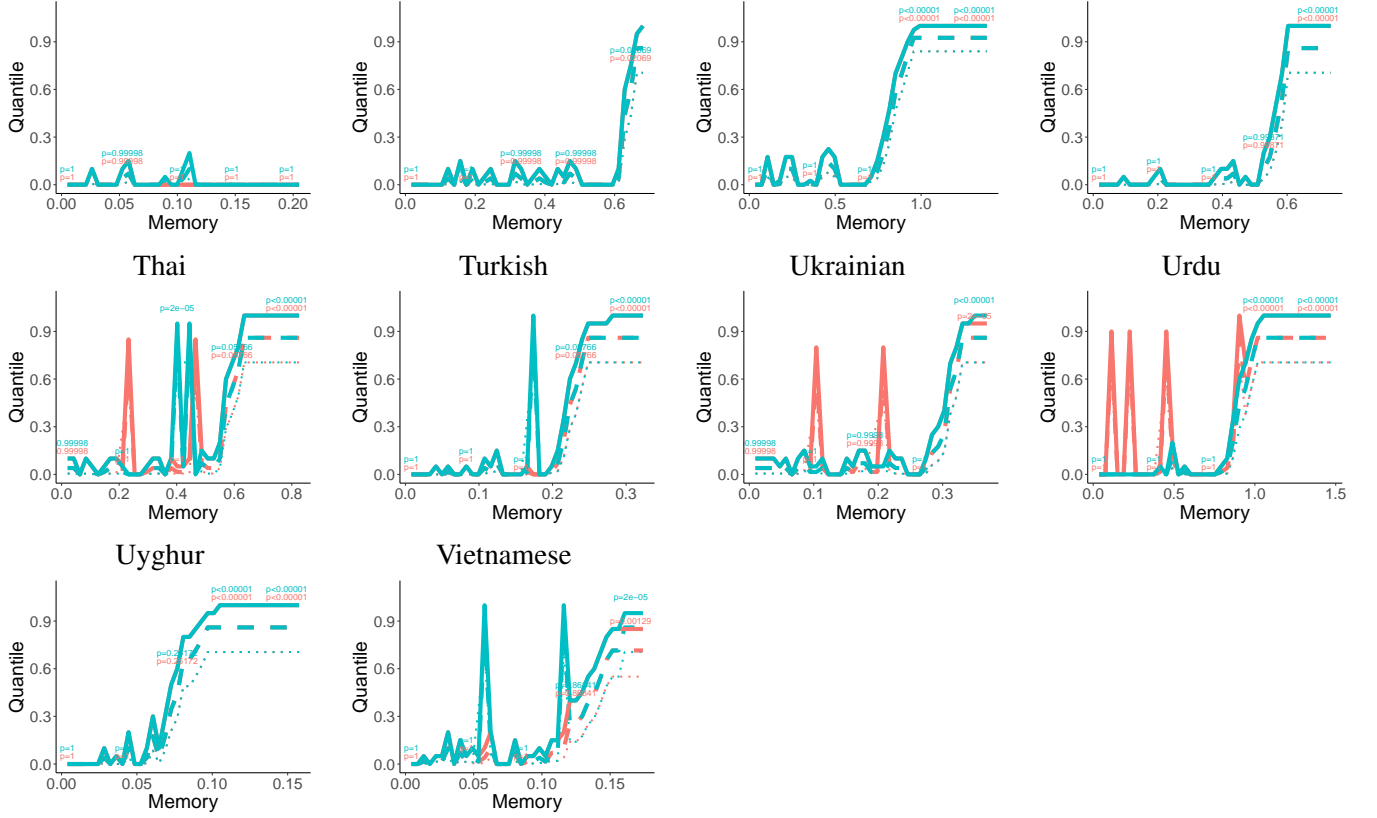


Figure 13: Quantiles: At a given memory budget, what percentage of the baselines results in higher listener surprisal than the real language? Solid curves represent sample means, dashed lines represent 95 % confidence bounds; dotted lines represent 99.9 % confidence bounds. At five evenly spaced memory levels, we provide a p-value for the null hypothesis that the actual population mean is 0.5 or less. Confidence bounds and p-values are obtained using an exact nonparametric method (see text).

{tab:quantil

## 9 Optimization

We cannot optimize for AUC using the method of (CITE), because we cannot construct an unbiased gradient estimator  $AUC \propto \sum_t (\sum_{s \leq t} I_s) t I_t = \sum_{s \leq t} t I_t I_s = I_1^2 + 2I_1 I_2 + 2I_2^2 + \dots$

As a surrogate, we propose to maximize  $I_1$ . This corresponds

It provides an accurate approximation to the AUC if  $I_s$  is small for  $s > 1$ , which holds for the n-gram based estimator.

If  $I_s < \epsilon$  for  $s > 1$ , then

$$I_1^2 \leq \sum_{1 \leq s \leq t \leq T} t I_t I_s \leq I_1^2 + I_1 \epsilon T^2 + \epsilon^2 T^2 \quad (50)$$

That means,  $I_1^2$  is an accurate approximation of the AUC when  $\epsilon T^2$  is small.

Softmax gradient update corresponds to adding  $\alpha$  to the target logit and removing  $\alpha$  from all other logits. Equivalently, add a certain dynamic amount to the logcount of the target and the total logcount.

Counting corresponds to adding 1 to the target probability.

## References

- Levy, R. (2008). Expectation-based syntactic comprehension. *Cognition*, 106(3):1126–1177.
- Lewis, R. L. and Vasishth, S. (2005). An activation-based model of sentence processing as skilled memory retrieval. *Cognitive Science*, 29(3):375–419.
- McElree, B. (2000). Sentence comprehension is mediated by content-addressable memory structures. *Journal of psycholinguistic research*, 29(2):111–123.
- Nicenboim, B. and Vasishth, S. (2018). Models of retrieval in sentence comprehension: A computational evaluation using bayesian hierarchical modeling. *Journal of Memory and Language*, 99:1–34.
- Vasishth, S., Nicenboim, B., Engelmann, F., and Burchert, F. (2019). Computational models of retrieval processes in sentence processing. *Trends in Cognitive Sciences*.