

Unpacking additivity
biases reveal the
class of sampling
algorithm for word
predictability

Yutong Zhang and Matthew Husband
University of Oxford

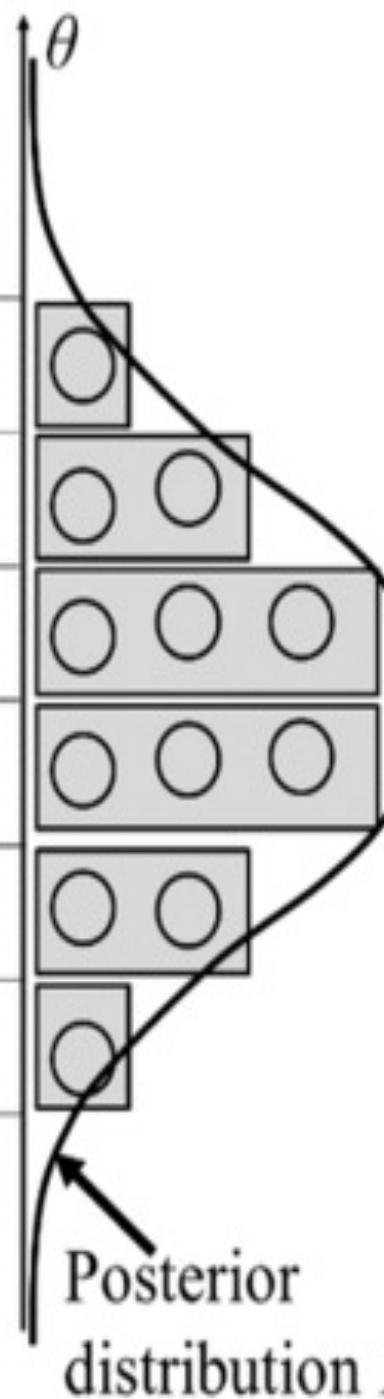


FACULTY OF
LINGUISTICS,
PHILOLOGY
AND
PHONETICS



Initial Sample (θ^0)

Prior distribution $p(\theta)$



Posterior
distribution

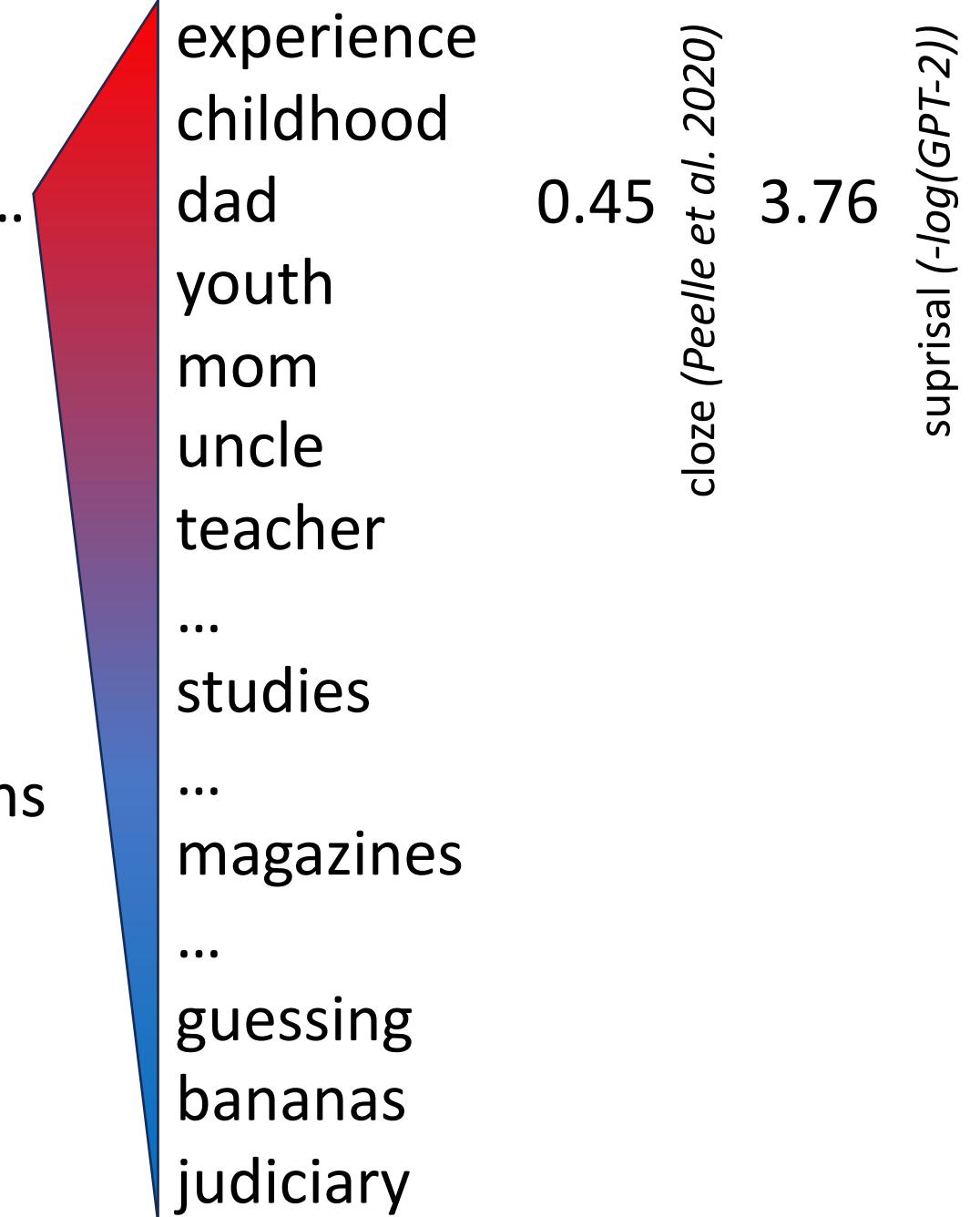
Jordan learned a lot about cars from his...

“dad” is *predictable* given this context.

- Cloze: 0.45 (Peelle et al 2020)
- GPT-2 surprisal: 3.76

As are many other potential continuations

Question: How do we arrive at ‘dad’?



Roadmap

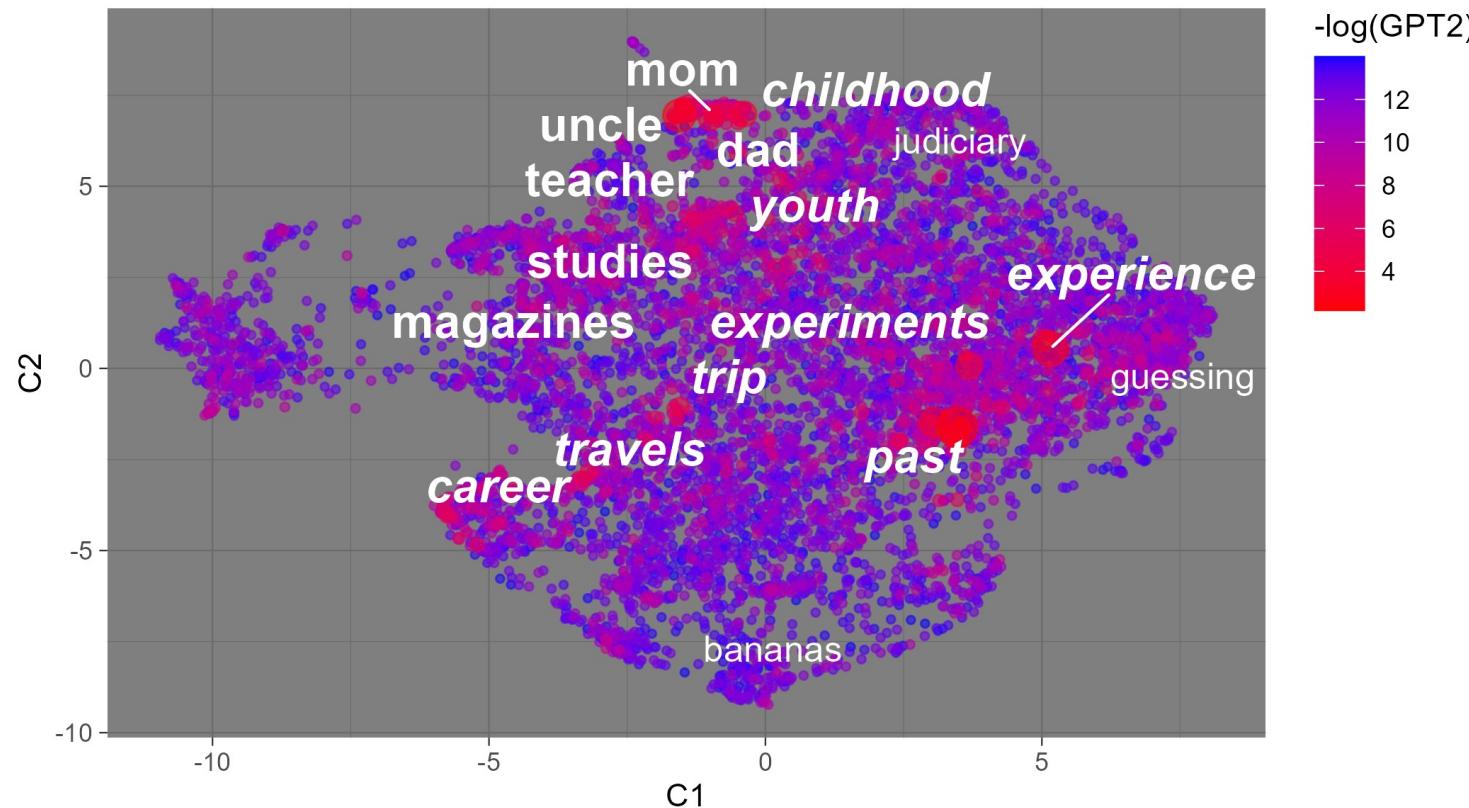
- Word Predictability as Bayesian Inference
- Classes of Sampling Algorithms for Approximate Bayesian Inference
- Biases in Human Judgement and Decision Making
- Initial Study examining Biases in Word Predictability
- Discussion & Conclusion

Word Predictability

- The difficulty of comprehending a word is related to its predictability in context.
- Computational accounts of predictability
 - Predictability reflects Bayesian updating of a prior probability distribution over lexical items, raising or lowering the probability of any lexical item given input (Kuperberg & Jaeger, 2016)
 - Likely a form of probabilistic inference, e.g. Surprisal Theory (Shain et al 2024)
- But such accounts face computational complexity (Kwisthout & van Rooij, 2020)
 - Can be computationally intractable to calculate Bayes' Rule.
 - Especially for large and complex space.

Word predictability – a complex space!

- Jordan learned a lot about cars from his...

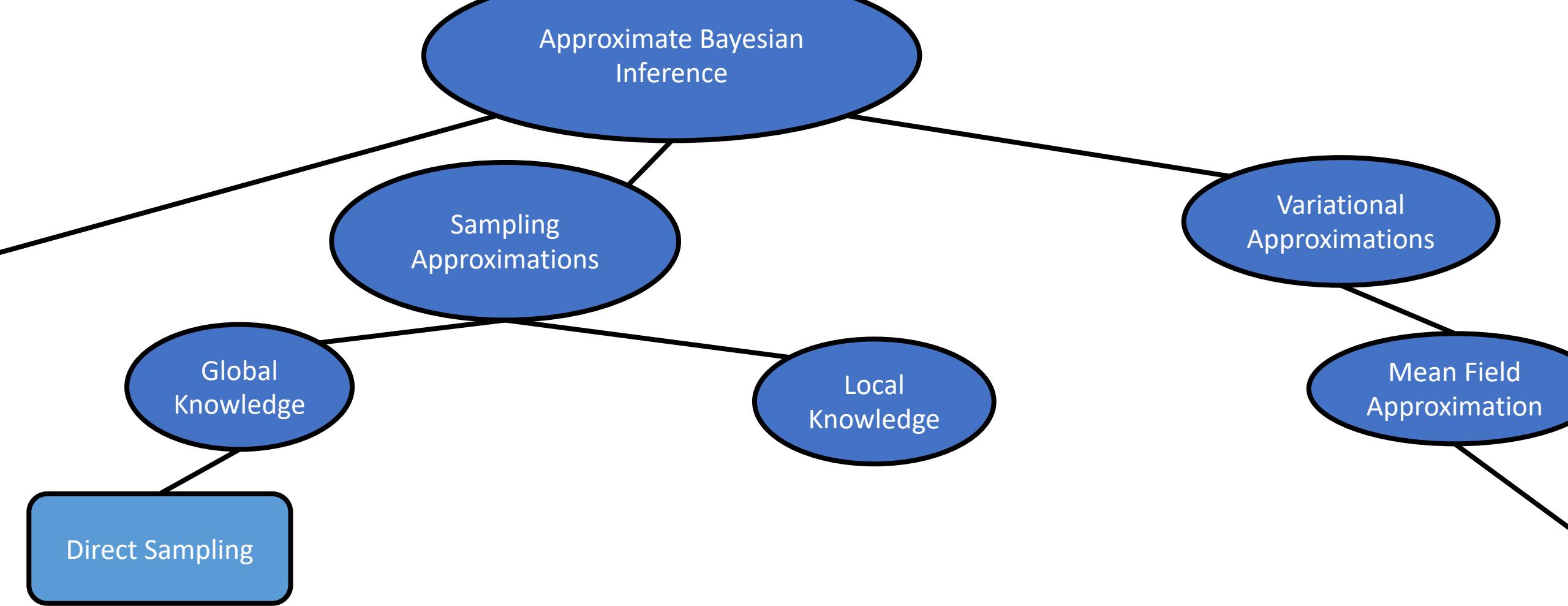


Word embedding: GLoVe 50
filtered for the 30,000 most
common English words, UMAP to
2 dimensions for visualization.

Surprisal calculated from GPT-2

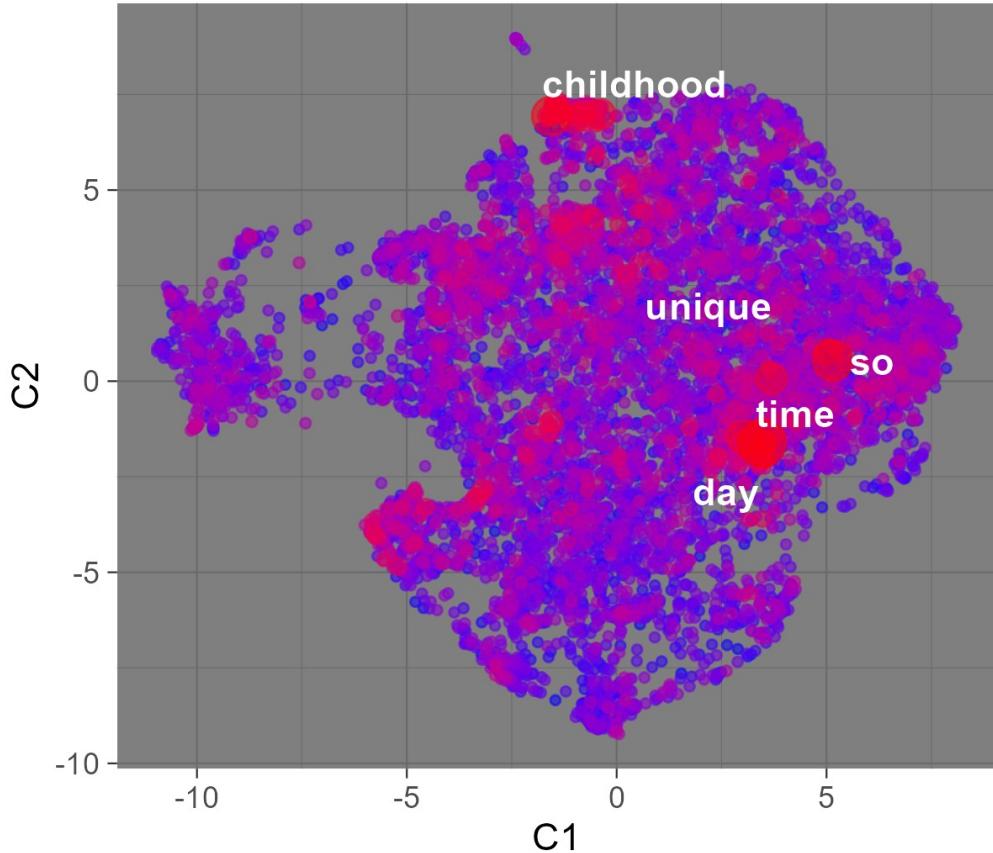
A Puzzle

- Comprehenders appear to be Bayesian about word predictability.
- BUT! Probabilistic (Bayesian) inference is likely to not be tractable in real time given the limits of our cognitive resources.
- The Question: How then does comprehension make use of probabilistic inference for word predictability if such inferences not tractable in real time under resource limitations?
 - Responding to this question seems likely to open up new avenues for research.

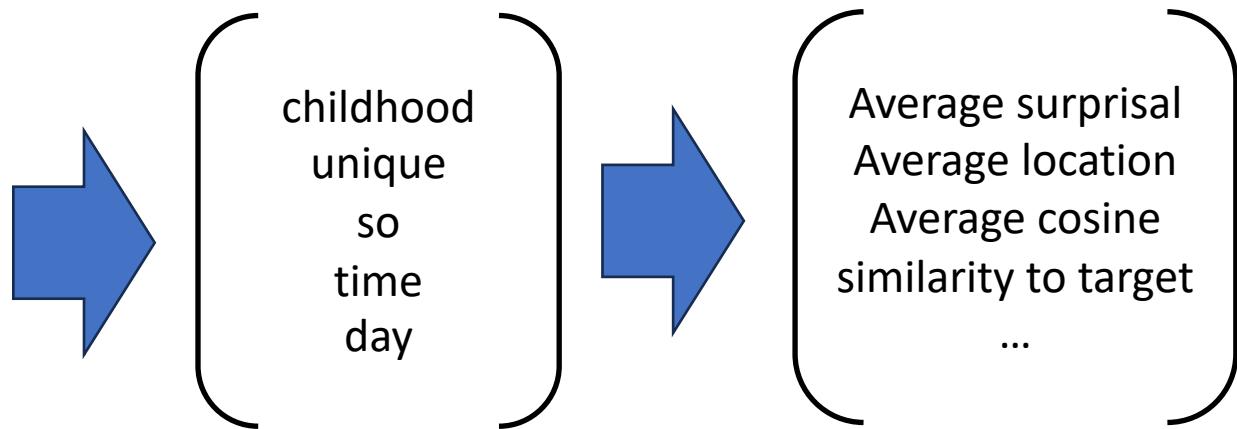


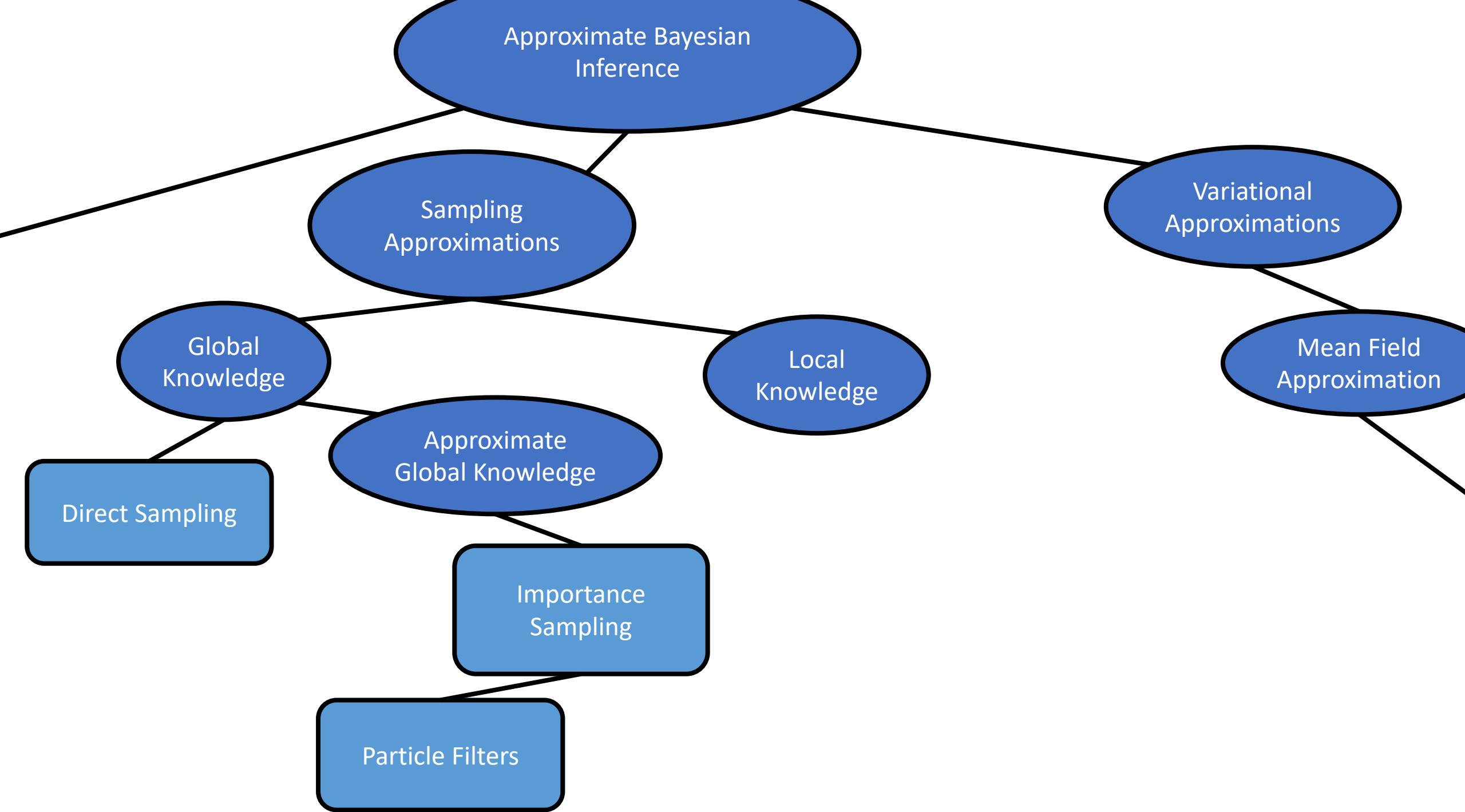
Direct Sampling

Jordan learned a lot about cars from his...



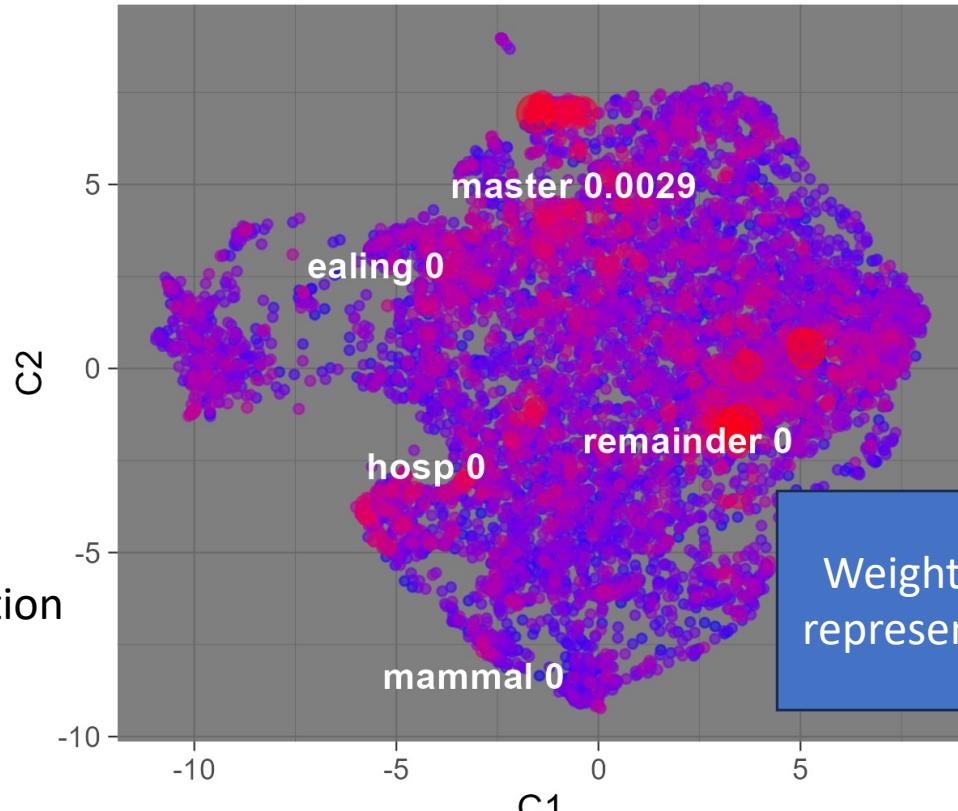
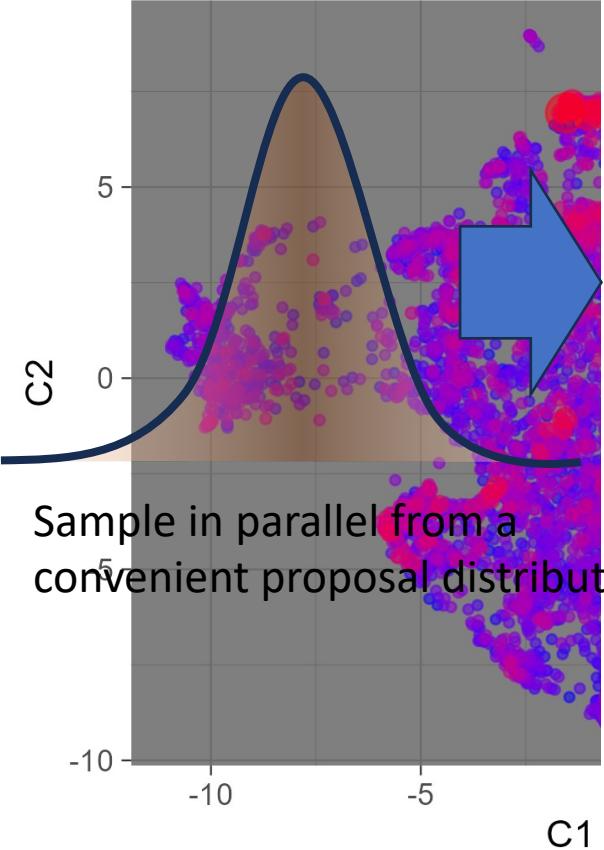
Sample in parallel from the posterior by probability





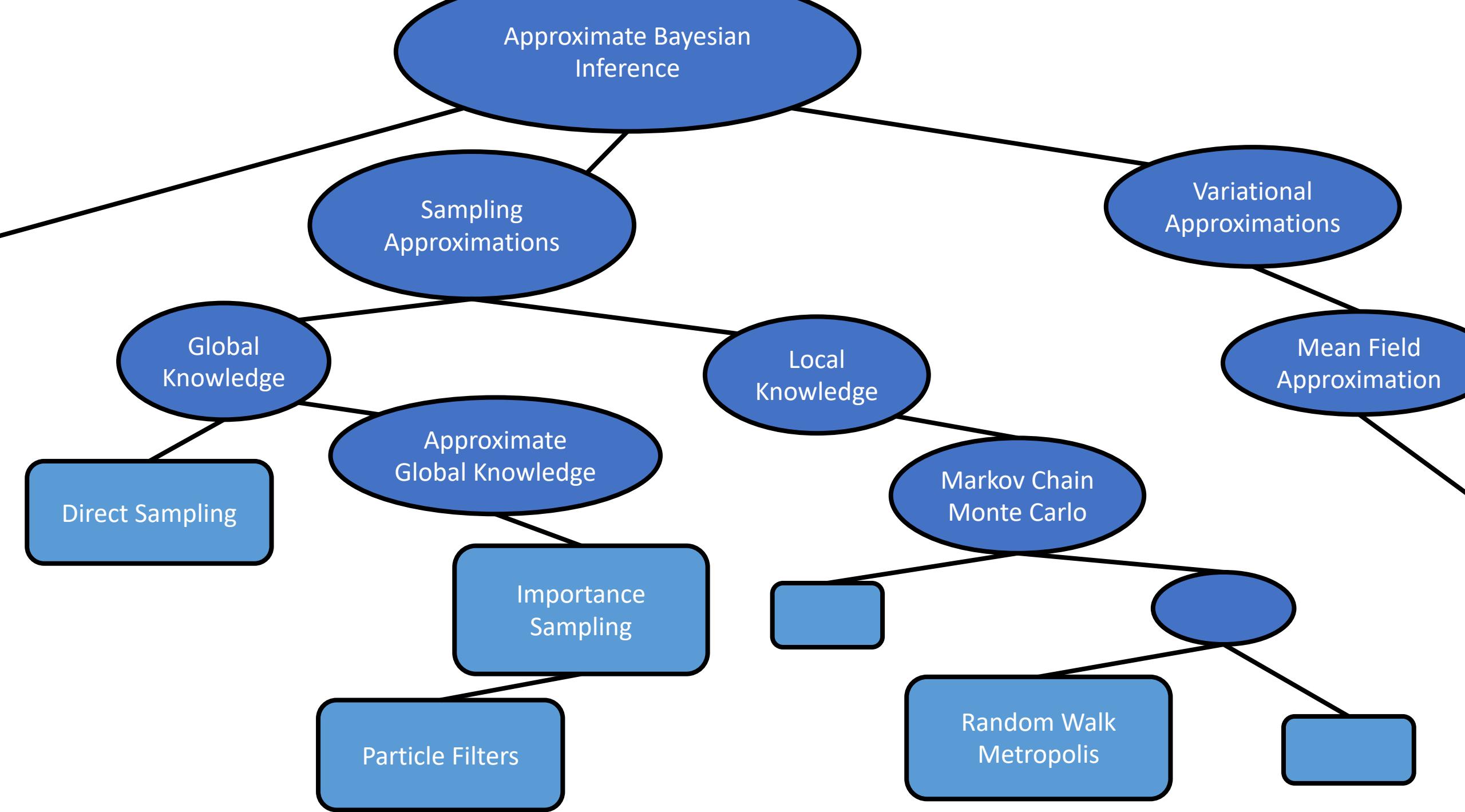
Importance Sampling

Jordan learned a lot about cars from his...



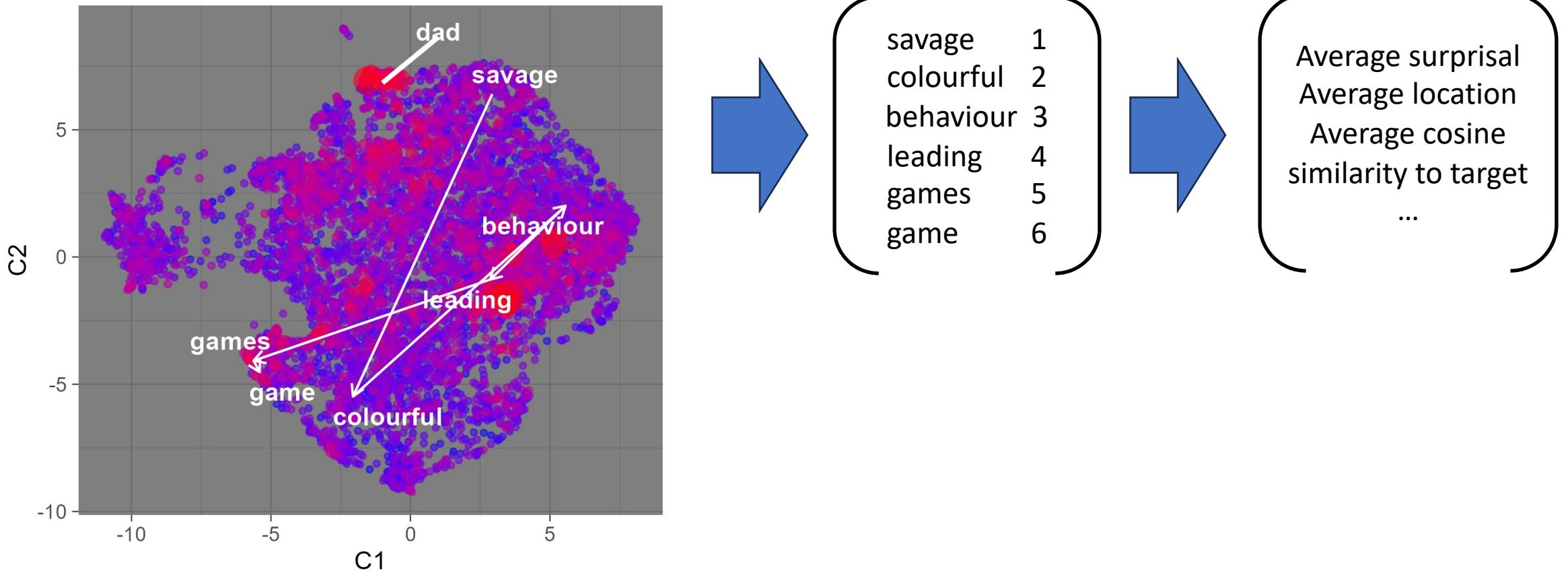
Weighting tends to over-represent typical examples

Average surprisal
Average location
Average cosine similarity to target
...



MCMC Sampling

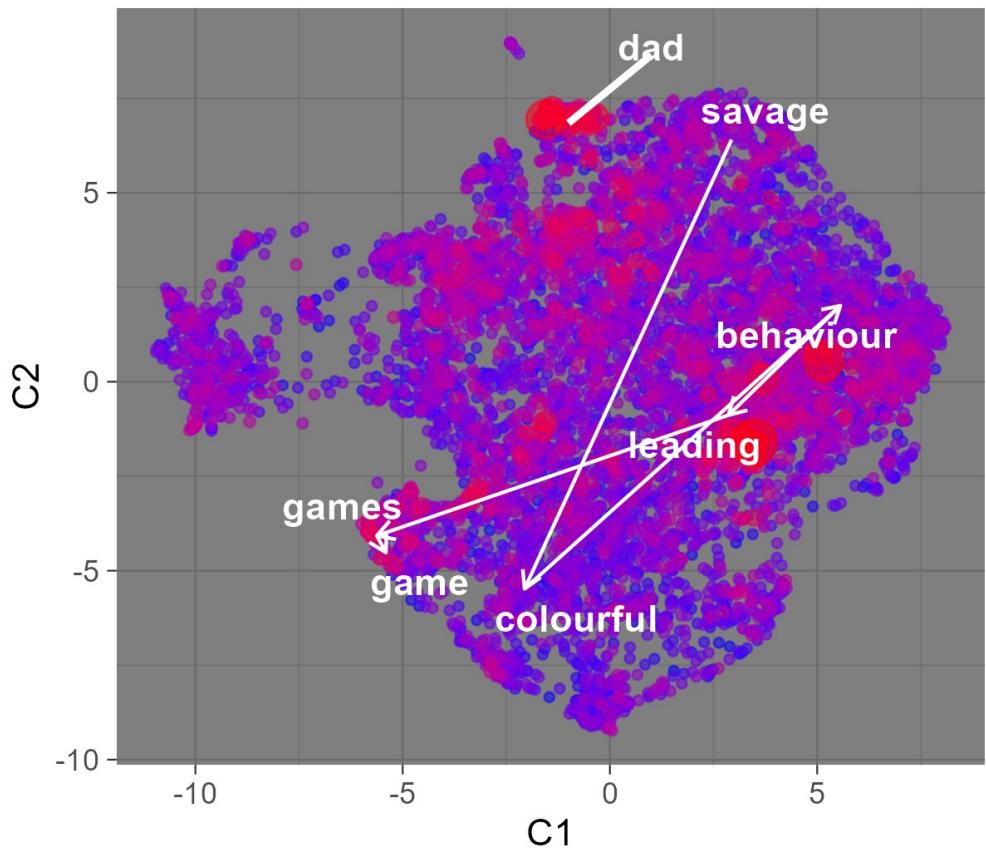
Jordan learned a lot about cars from his...



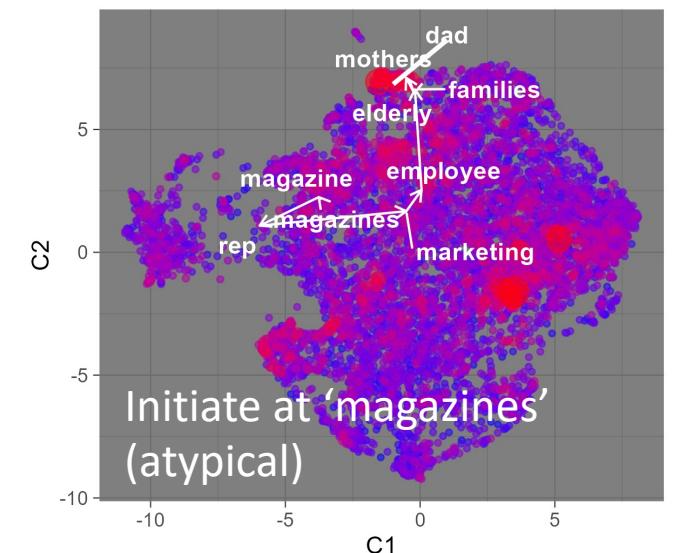
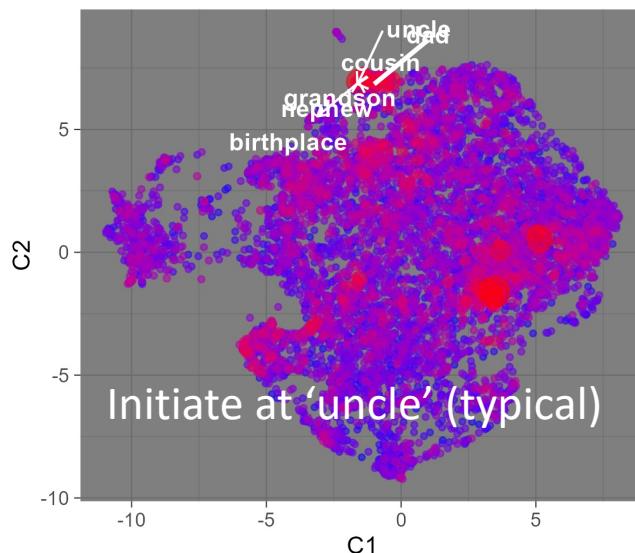
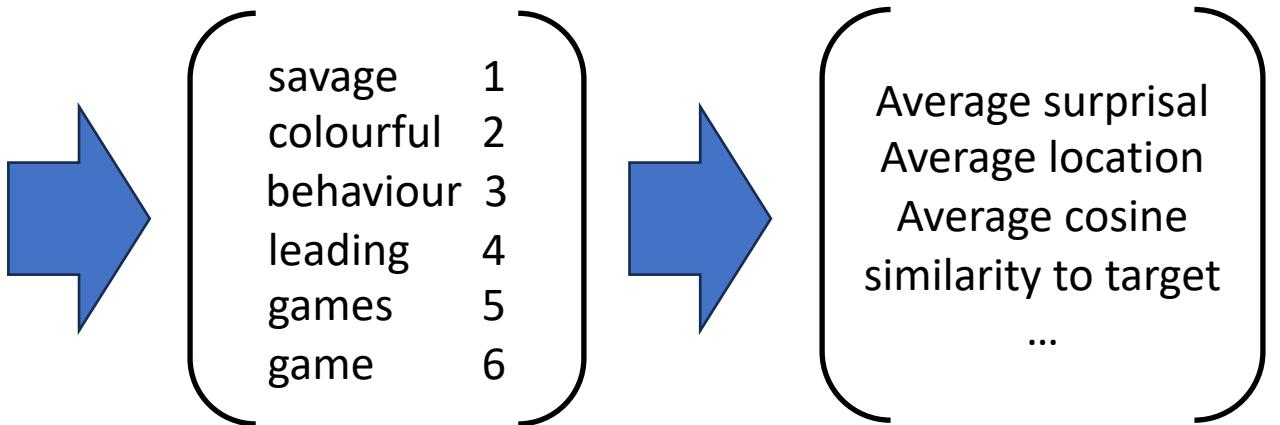
Initiate a chain with a sample. Propose a new sample;
accept/reject that sample given its relative probability
to the current sample in the chain

MCMC Sampling

Jordan learned a lot about cars from his...



Initiate a chain with a sample. Propose a new sample; accept/reject that sample given its relative probability to the current sample in the chain



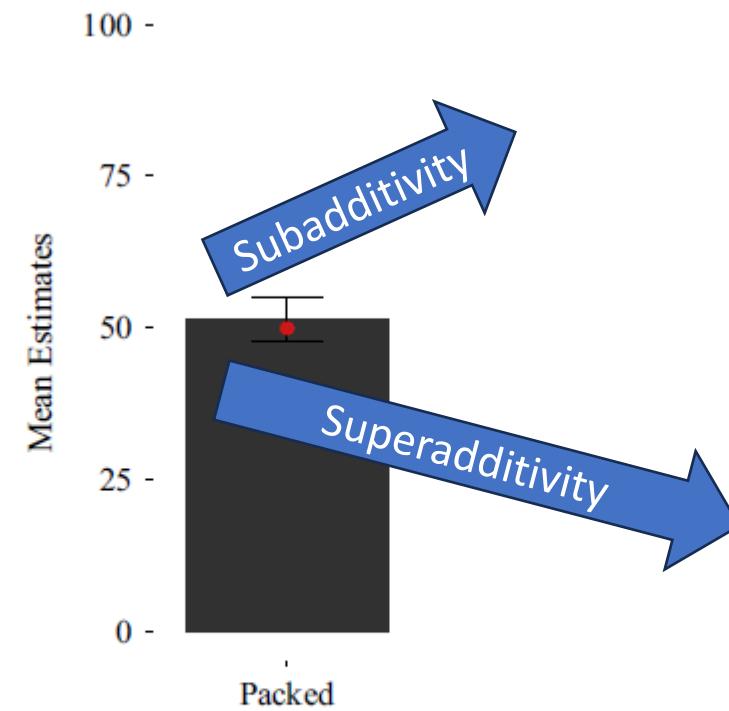
Sampling algorithm	Knowledge	Sytematic Bias	Areas of explanation
Direct Sampling	Global	Stochastic behavior	Probability matching behavior (Gershman, 2014); Exploration-exploitation trade-off (Gershman, 2018; Speekenbrink & Konstantinidis, 2015)
Importance Sampling	Approx. Global	Subadditivity (Over-represents typical samples)	Perceptual stimuli reproduction (Shi et al., 2010); Semantic memory dependence and working memory load (Lloyd et al., 2019; Sanborn et al., 2010)
MCMC Sampling	Local	Subadditivity and Superadditivity (Over-represents /under-represents typical examples when initiated with a typical/atypical example)	Anchoring bias (Lieder et al., 2018); Biases in probability judgments (Dasgupta et al., 2017; Sanborn & Chater, 2016);

Question: Are there systematic biases in word predictability, and can we use it to narrow down the class of sampling algorithm that supports it?

Outside of language comprehension, the probabilistic inference computations underlying judgment and decision making in human cognition are known to be systematically biased.

Probability judgments and unpacking

- I see a table.
- What is the probability that I also see...
 - an object starting with the letter C? [packed]
 - a chair, computer, curtain, or any other object starting with the letter C? [unpacked – typical]
 - a cannon, cow, canoe, or any other object starting with the letter C? [unpacked – atypical]



Summary

- Unpacking context impacts probability judgments
 - Judgments are *subadditive* when unpacked examples are typical, leading to an increase in overall probability
(Dasgupta et al 2017; Hadjchristidis et al 1999; Fox & Tversky 1998; Tversky & Koehler 1994).
 - Judgments are *superadditive* when unpacked examples are atypical, leading to a decrease in the overall probability judgment
(Dasgupta et al 2017; Dougherty & Hunter 2003 ; Sloman et al 2004).
- Human Probability Judgment shows systematic biases.
 - These are proposed to arise from Approximate Probabilistic Inference algorithms e.g. consistent with the MCMC class (Dasgupta et al 2017).

Back to word predictability

- If comprehension is also a form of Approximate Probabilistic Inference, does it show biases similar to those found in judgment and decision making?
- If so, do those biases diagnose which class of sampling algorithms is a more appropriate fit for word predictability?

Methods

- 60 Participants; 120 Items; 3 conditions
- Range of target word cloze: 0.30 to 0.70 (Peelle et al. 2020)
- Unpacking: Typical unpacking was a high cloze word; atypical unpacking was a low cloze word

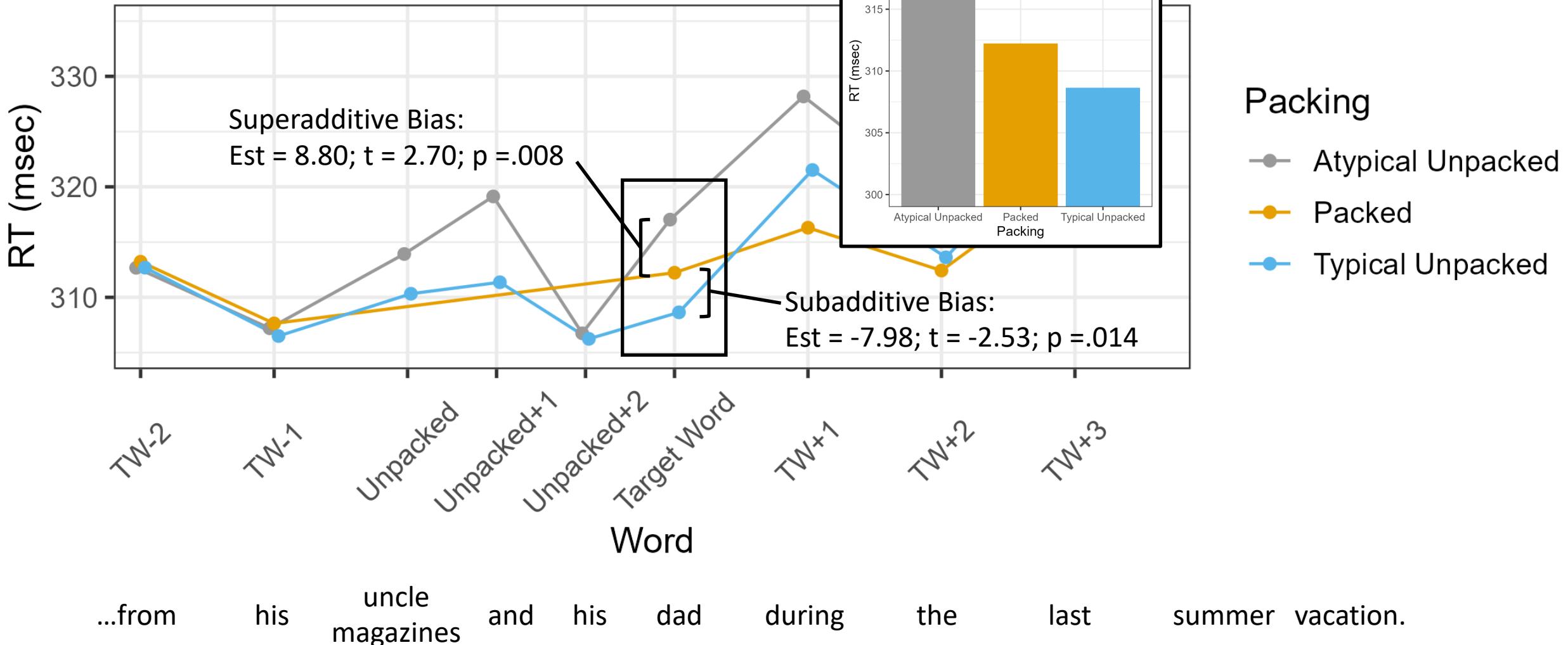
Cloze	Context	Target Word
0.30	The boy never accomplished anything since he lacked the	<u>motivation</u> to pursue his goals
0.51	Her friend said the story was nothing but	<u>lies</u> that could be seen through easily.
0.70	The bride spent weeks shopping for her	<u>dress</u> for the upcoming wedding.

- Also measured Word2Vec and LSA similarity between the unpacked word and target word.

Predictions

Bias	Description	Reading Time Prediction
Subadditivity	Perceived predictability of a word is <i>higher</i> when the context is unpacked to <i>typical</i> examples.	Faster than packed baseline
Superadditivity	Perceived predictability of a word is <i>lower</i> when the context is unpacked to <i>atypical</i> examples.	Slower than packed baseline

Jordan learned a lot about cars...



Is this just priming?

- LSA between unpacked word and target word
 - Inclusion of centered LSA scores did not interact with Typicality
 - or reduce the effect of Typicality significantly.
- Word2Vec between unpacked word and target word
 - Inclusion of centered word2vec scores reduced the significant of Typicality,
 - but did not significantly interact with Typicality and did not improve model fit ($\chi^2 = 1.203$; $p = .273$).
 - Typicality was still significant with word2vec as a covariate

	Est	t	p
Typicality	8.38	3.07	.003
LSA	0.78	0.06	.951
Typicality:LSA	8.83	0.42	.674

	Est	t	p
Typicality	5.14	0.89	.375
w2v	16.70	1.53	.128
Typicality:w2v	23.01	1.09	.275

	Est	t	p
Typicality	10.40	3.23	.001
w2v	10.97	1.14	.254

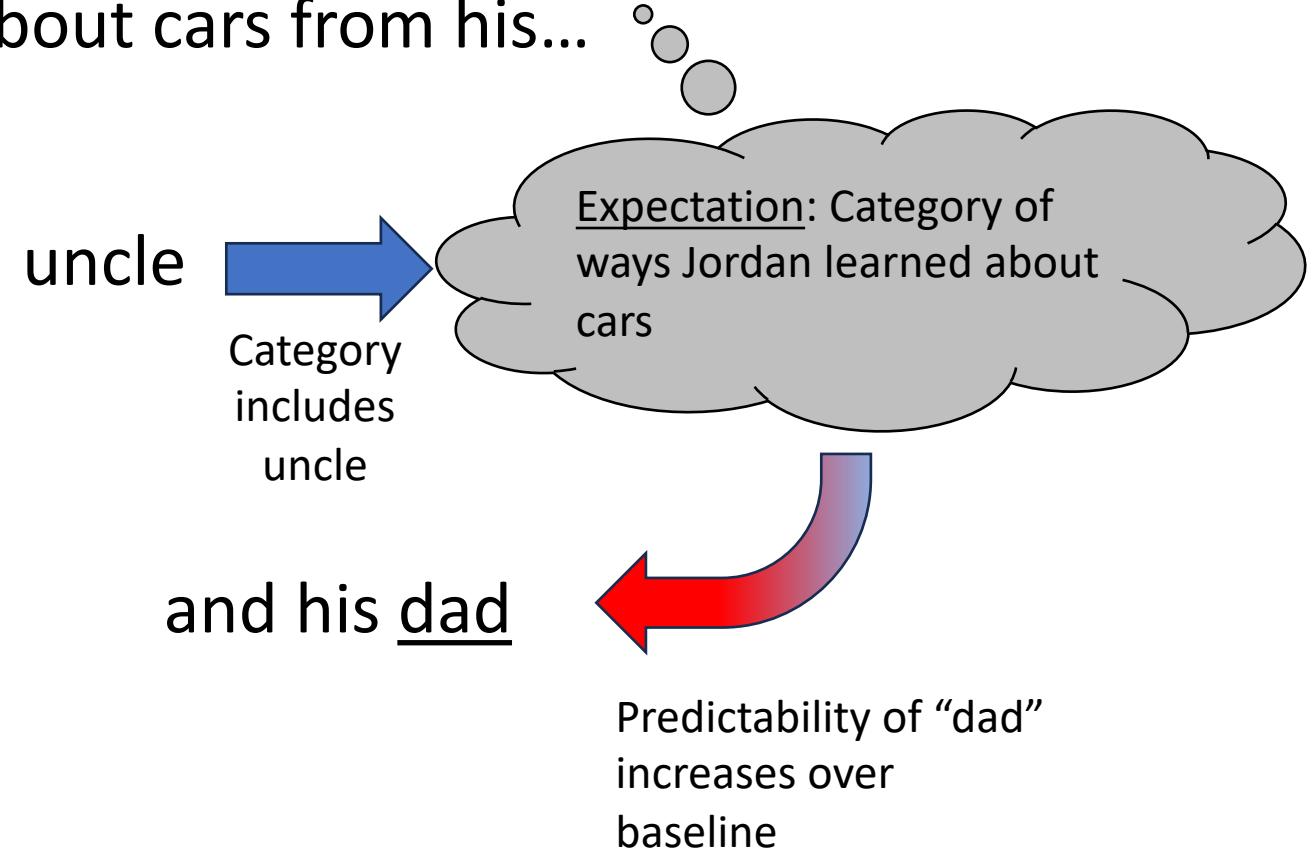
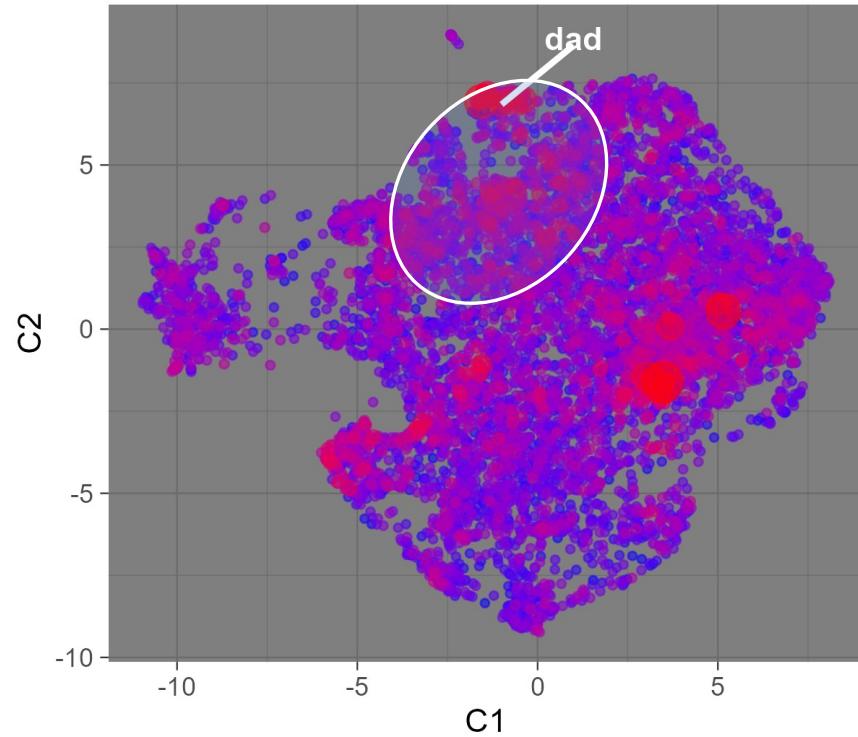
What about surprisal?

- Surprisal of the target word in each condition was measured with GPT-2.
- Surprisal did not interact with Unpacking Typicality
 - and did not improve model fit ($\chi^2 = 2.760$; $p = .252$)
 - Typicality remained significant with GPT-2 Surprisal as a covariate.

	Est	t	p
Typical	-9.46	-2.66	.009
Atypical	9.28	2.60	.011
GPT2	-4.63	-0.15	.882
Typical:GPT2	86.77	1.05	.300
Atypical:GPT2	-64.86	-0.63	.533

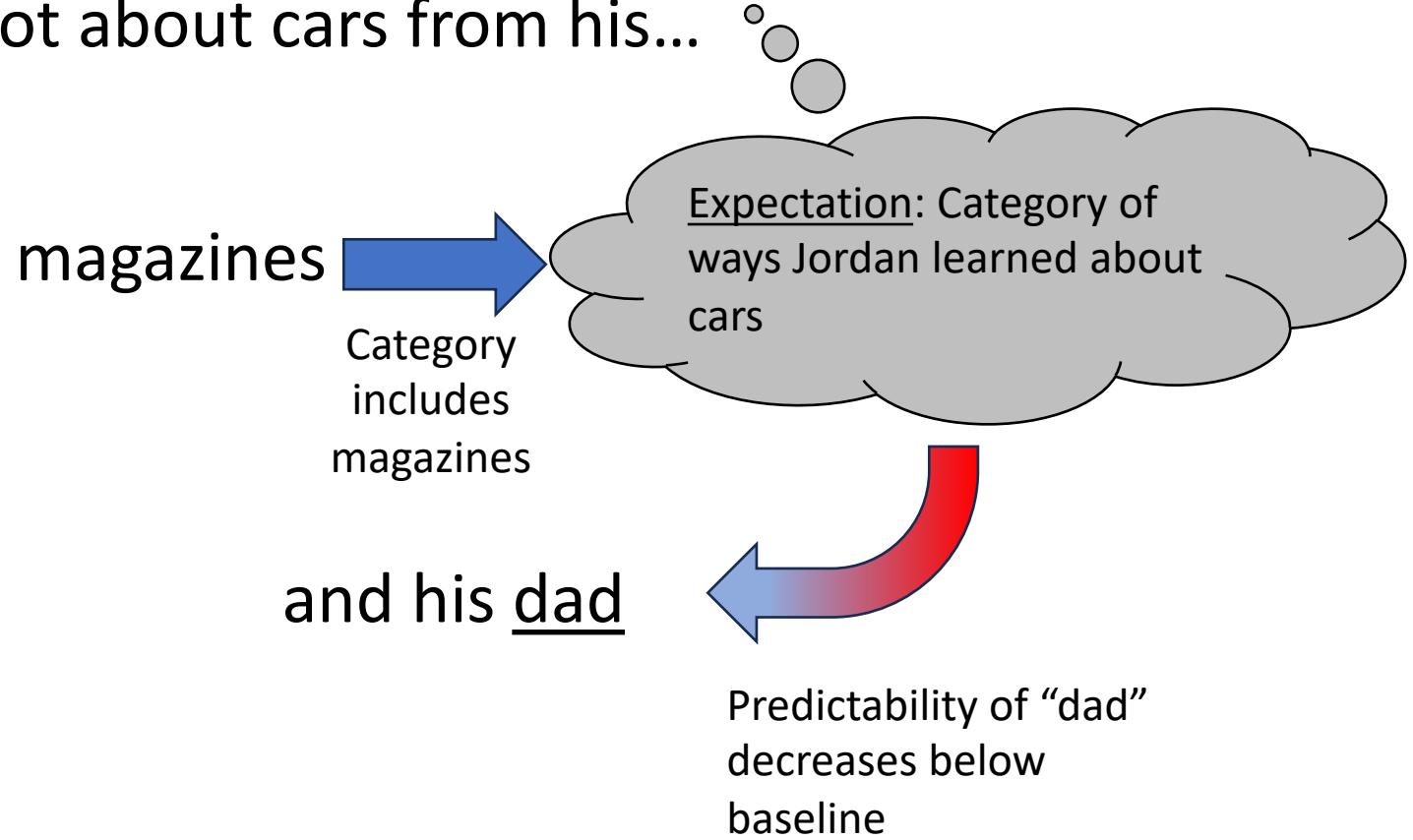
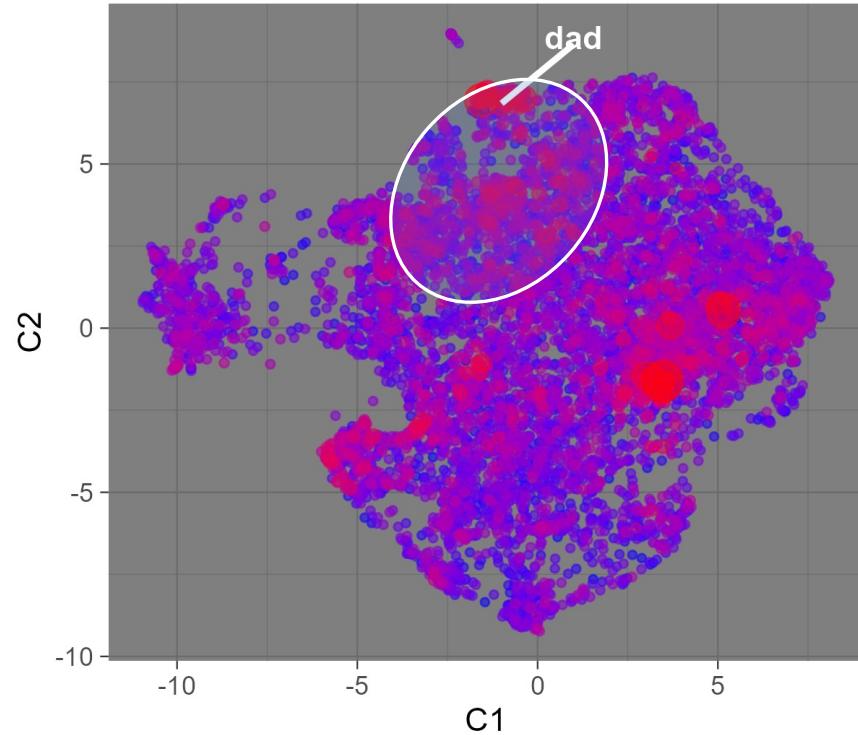
Unpacked Typical

Jordan learned a lot about cars from his...



Unpacked Atypical

Jordan learned a lot about cars from his...



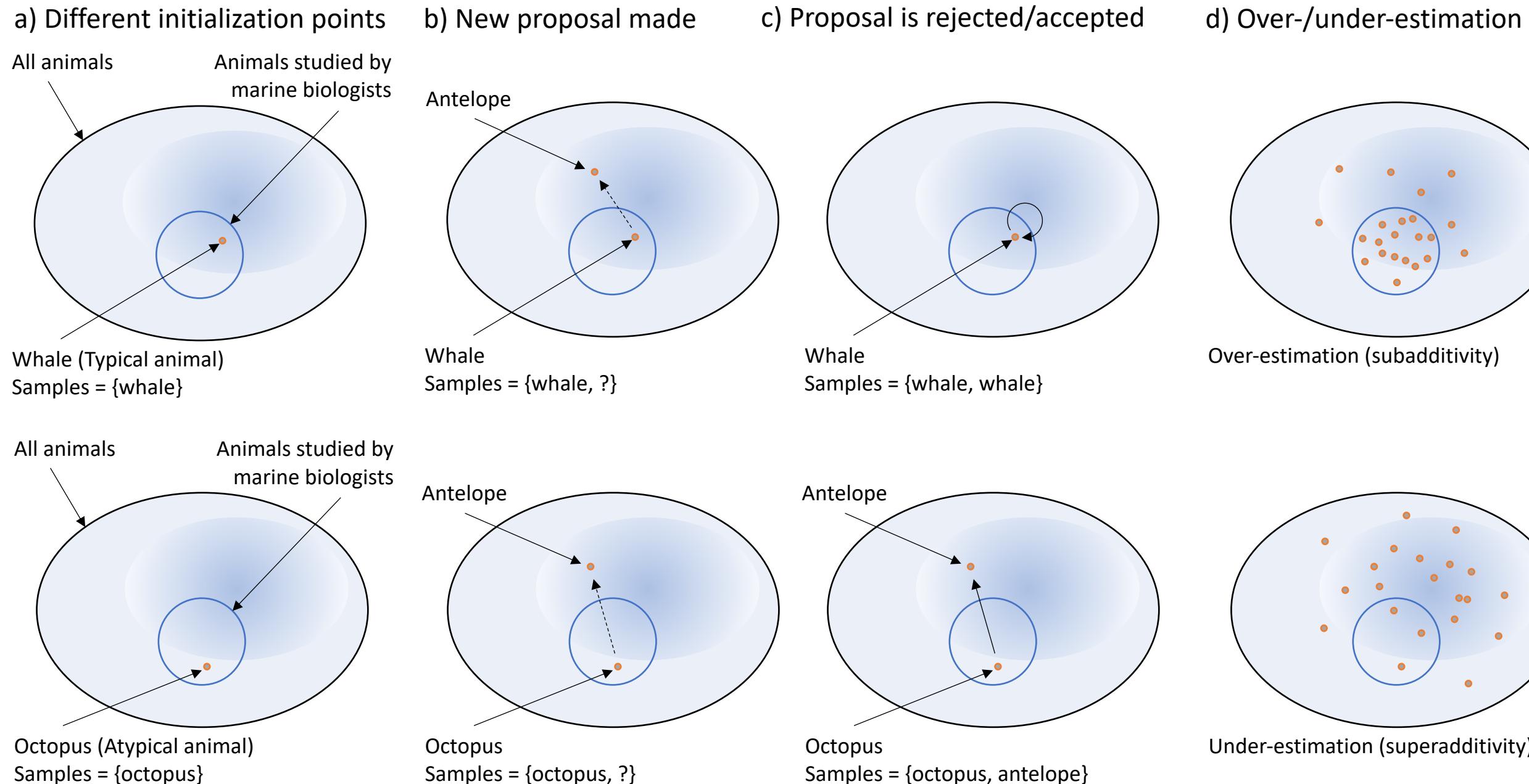
Word predictability by MCMC?

Why should language work this way?

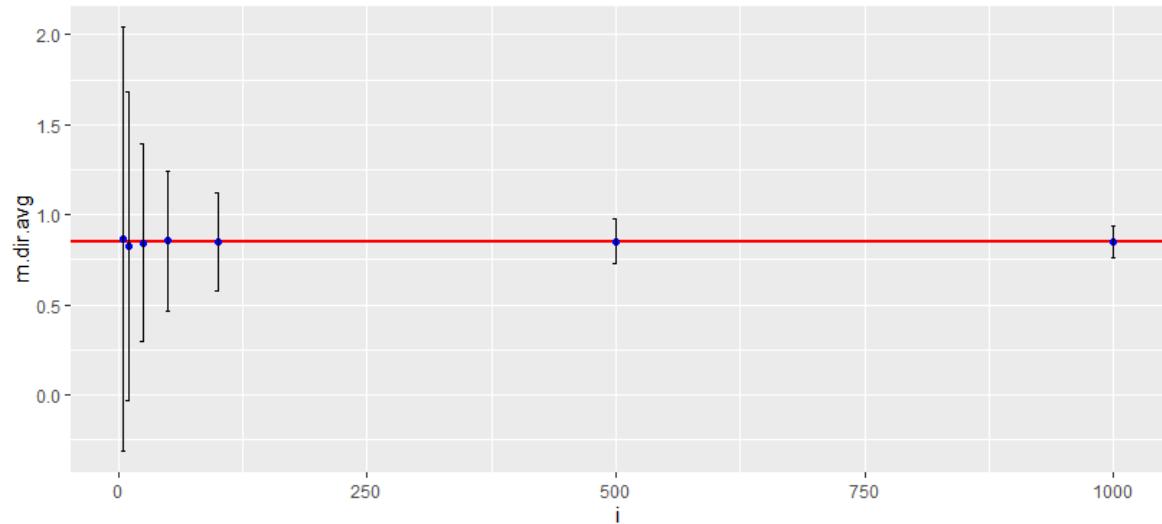
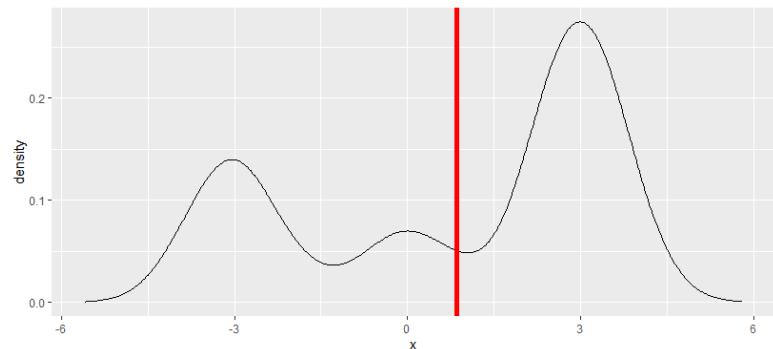
- Local knowledge
 - Only the relative probabilities of the current sample and the proposal are needed (Sanborn & Chater 2016; Stewart et al 2006)
- Makes appropriate use of prior context to directly anchor a sample chain
 - Such chains might even be reused for further computation down stream (Dasgupta et al 2018)
- Computational efficient
 - MCMC sounds fancy, but it has a straightforward implementation and is even biologically plausible (Buesing et al 2011; Moreno-Bote et al., 2011; Pecevski et al 2011)

Conclusions

- Comprehenders may be Bayesian about word predictability.
 - But they are unlikely to be able to fully execute the computations involved due to resource limitations.
- Probabilistic inference can be approximated by a variety of algorithms
 - These algorithms have certain biases when resource limited.
- What biases do we see during human language comprehension?
 - Some evidence today that word predictability shows both subadditive and superadditive biases.
 - To the extent that sampling algorithms are appropriate models for cognition, Markov Chain Monte Carlo (MCMC) Sampling provides a better fit than Direct or Importance Sampling.



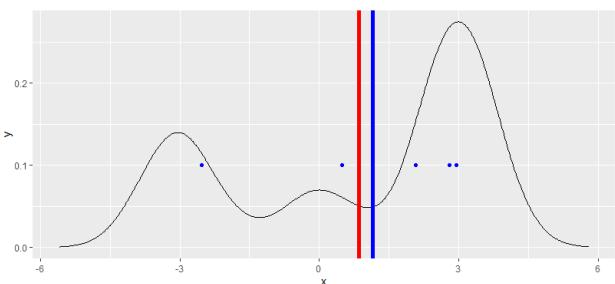
Direct Sampling



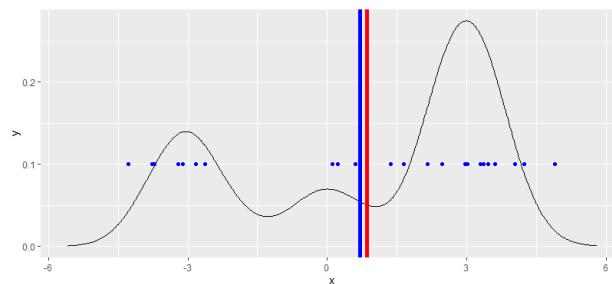
Take samples directly from the posterior distribution.

- Need to have global knowledge of the posterior
- No systematic bias with small numbers of samples – just added noise.

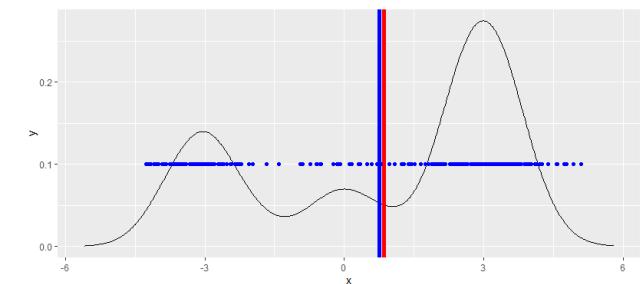
$N=5$



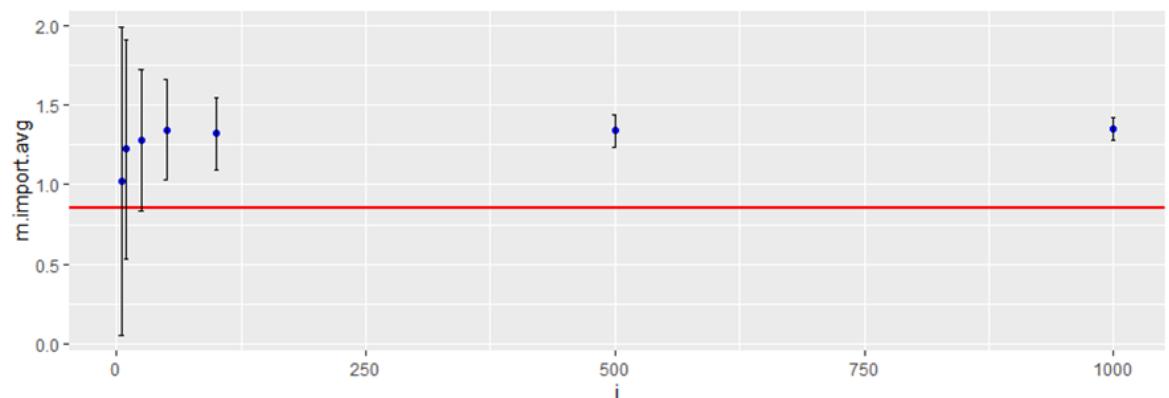
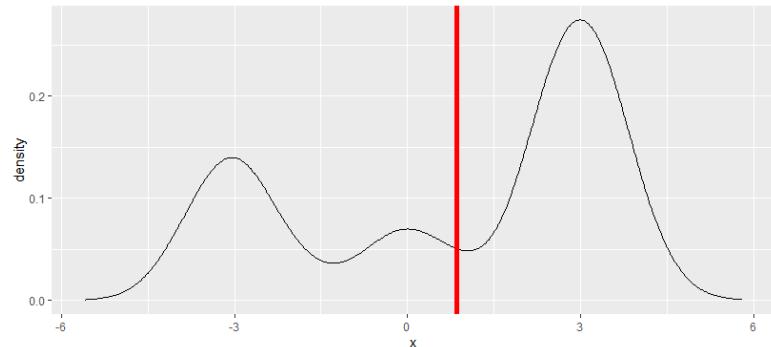
$N=25$



$N=250$



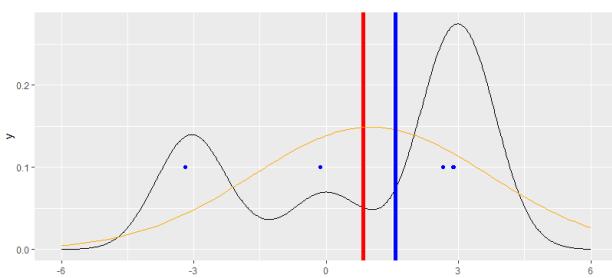
Importance Sampling



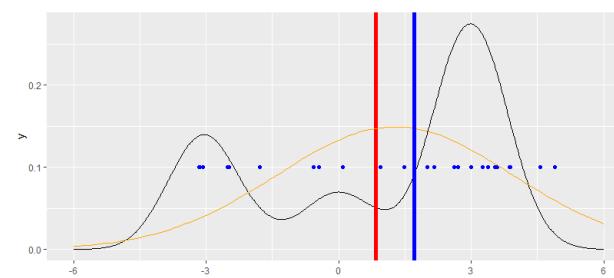
Take samples from the a proposal distribution that is (hopefully) related to the posterior distribution.

- Need to have approximate global knowledge of the posterior
- *Subadditivity*: Can be biased by more typical values.

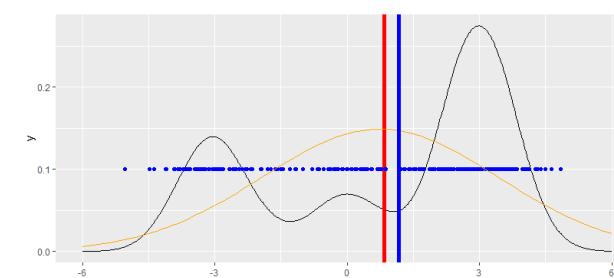
$N=5$



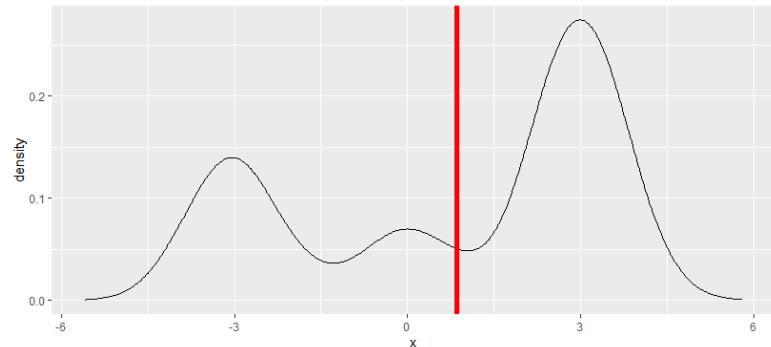
$N=25$



$N=250$



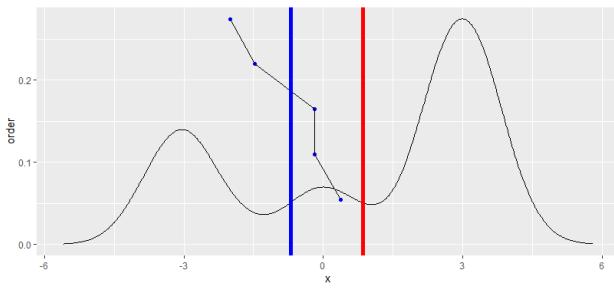
MCMC Sampling



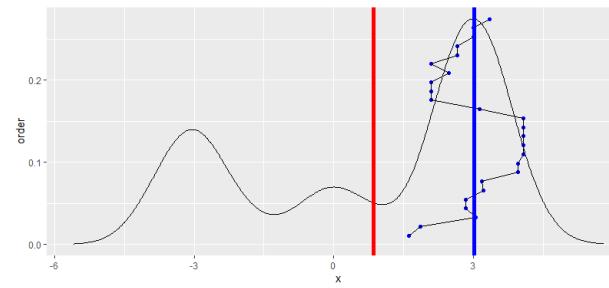
Initiate with one sample, propose a second, and compare their probabilities.

- Need to only have local knowledge of the posterior.
- *Subadditivity*: Can be bias by starting at a typical value.
- *Superadditivity*: Can be bias by starting at a more atypical value.

$N=5$



$N=25$



$N=250$

