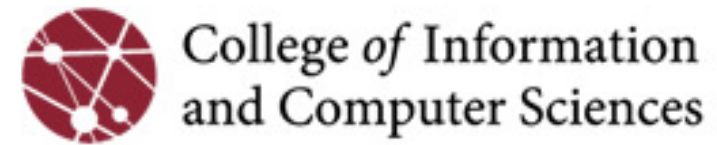




UMassAmherst



Commonsense Knowledge Representation in NLP



Xiang Lorraine Li

PhD Student at UMass Amherst,

Advised by Prof Andrew McCallum

Commonsense: Essential in AI ^[1]



A lake divides into two lakes when the water level falls

A five-year-old would know this!



A five year old would also know the brick will fall if it's not placed correctly

Common sense is sound practical judgement

- **Concerning everyday matters**
- **Basic ability to perceive, understand, and judge**
- **Shared by ("common to") nearly all people. ----Wikipedia**

Disclaimer: we mean, roughly, what a typical five year old knows about the world, including fundamental categories like time and space, and specific domains such as physical objects and substances; plants, animals, and other natural entities; humans, their psychology, and their interactions; and society at large. We will not attempt to be precise about this, but let us indicate roughly which issues we are considering and which we are ignoring. Obviously, this body of knowledge in fact depends on place, time, culture, social standing, personal characteristics (e.g. unusual cognitive or physical abilities or disabilities), schooling, perhaps on language. We ignore all that; without embarrassment, we have in mind a 21st-century, first-world, urban, with an appropriate level of schooling.

[1] <https://cacm.acm.org/magazines/2015/9/191169-commonsense-reasoning-and-commonsense-knowledge-in-artificial-intelligence/fulltext?mobile=false#F1>

Commonsense: Essential in AI [1]

- Natural Language Processing

- Machine T



Danish Pruthi @danish037 · 11月14日

Not being able to find my phone, I ask for help.

Me: Hey Google, could you call me, please?

Nest mini: you sure?

Me: Yes!

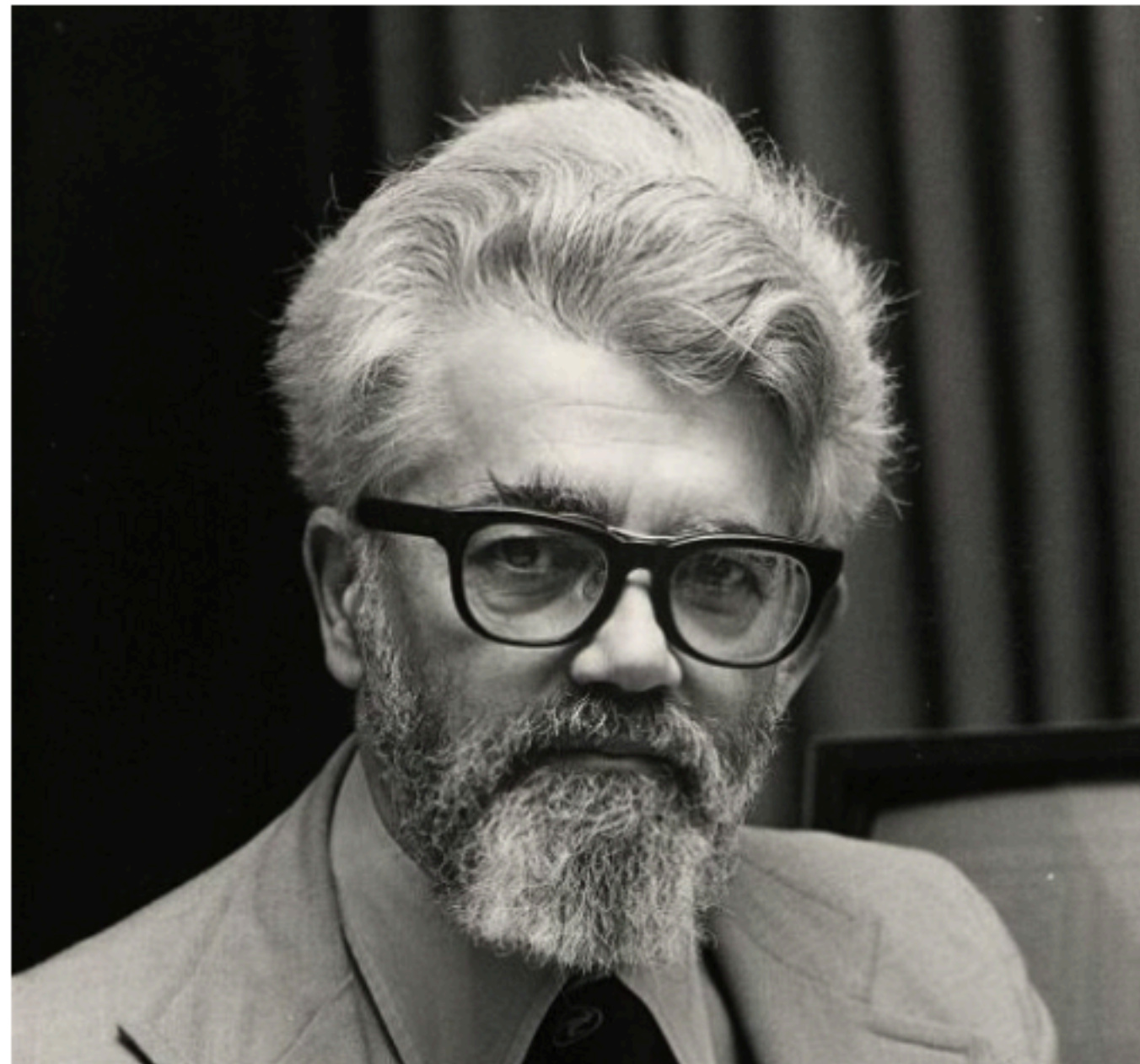
Nest mini: okay, I will call you please from now on.



- Smart Ho



Logical Formalization of Commonsense



John McCarthy

Advise Taker

1. First, we have a predicate "at". " $at(x, y)$ " is a formalization of " x is at y ". Under this heading we have the premises

$at(I, desk)$ (1)

$at(desk, home)$ (2)

$at(car, home)$ (3)

$at(home, county)$ (4)

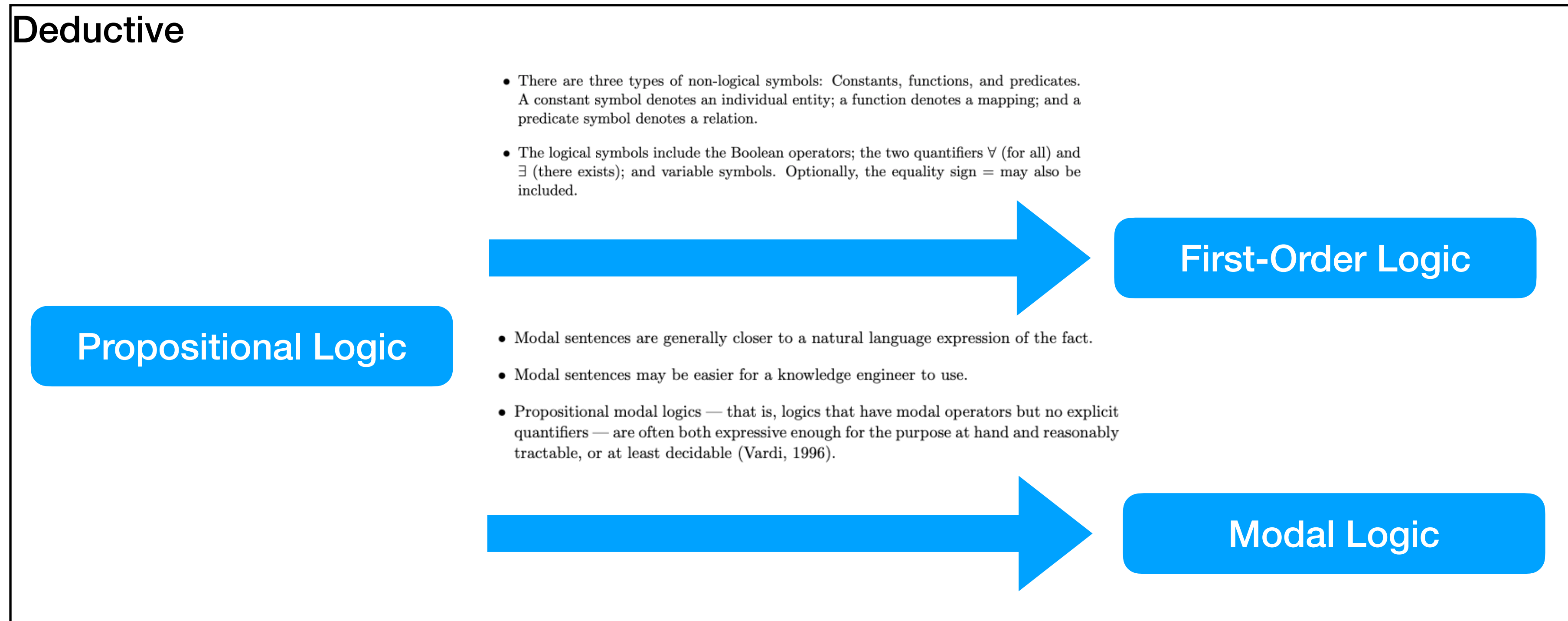
$at(airport, county)$ (5)

In order for a program to be capable of learning something it must first be capable of being told it — John McCarthy [1]

[Programs with Common Sense](#) was probably the first paper on logical AI, i.e. AI in which logic is the method of representing information in computer memory and not just the subject matter of the program. The paper was given in the Teddington Conference on the Mechanization of Thought Processes in December 1958 and printed in the proceedings of that conference. **It may also be the first paper to propose common sense reasoning ability as the key to AI.**

Logical Formalization of Commonsense

Incomplete list of different logic systems.



Plausible Reasoning

Non-monotonic Logic

Probabilistic Logic

Fuzzy Logic

Logical Formalization of Commonsense

But ... it did not end up well.

Past failures (in 70s – 80s) are inconclusive

- weak computing power
- not much data
- not as strong computational models
- not ideal conceptualization / representations

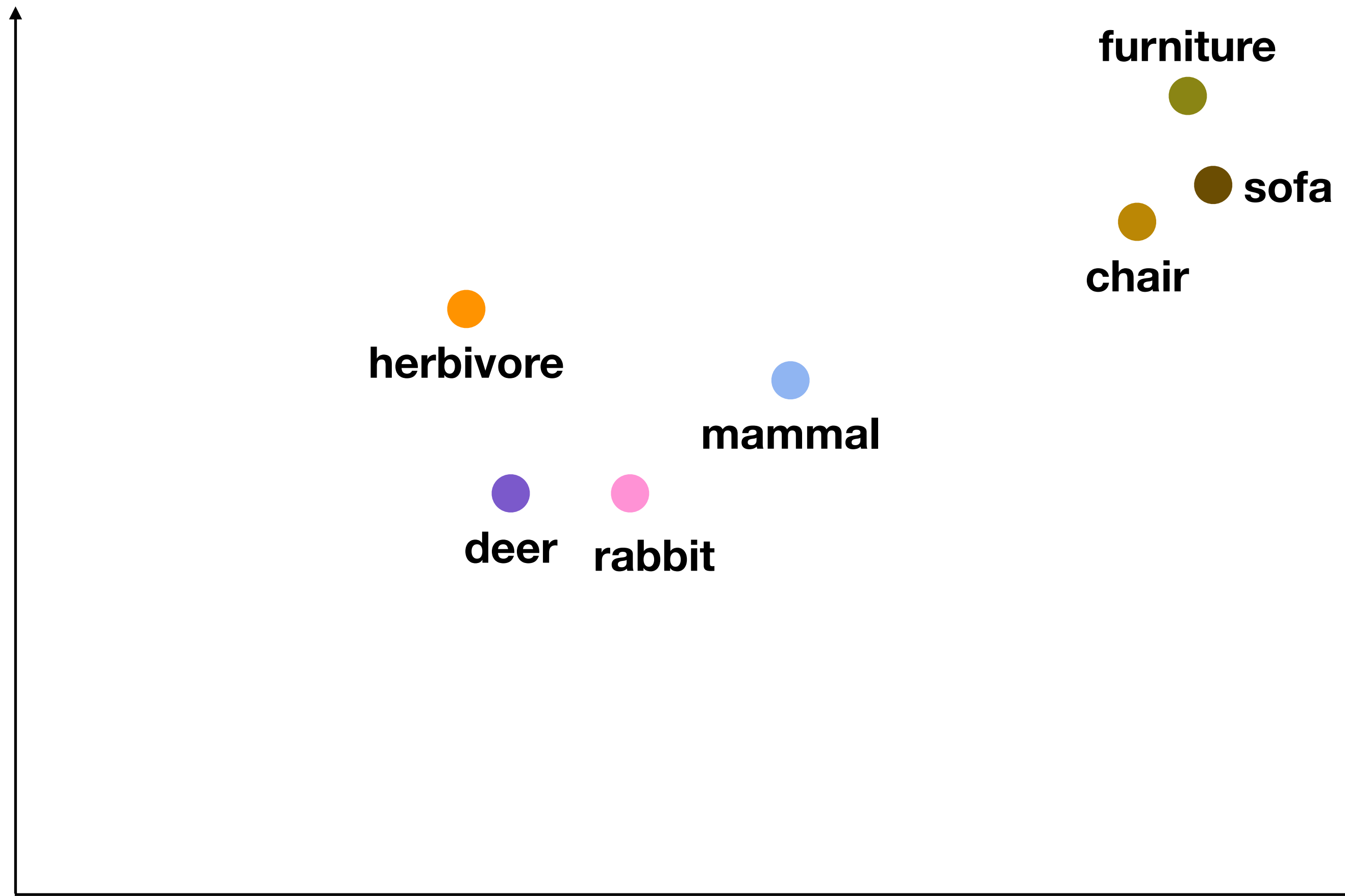


“I was told not to speak the word commonsense...” —Yejin Choi [1]

Outline

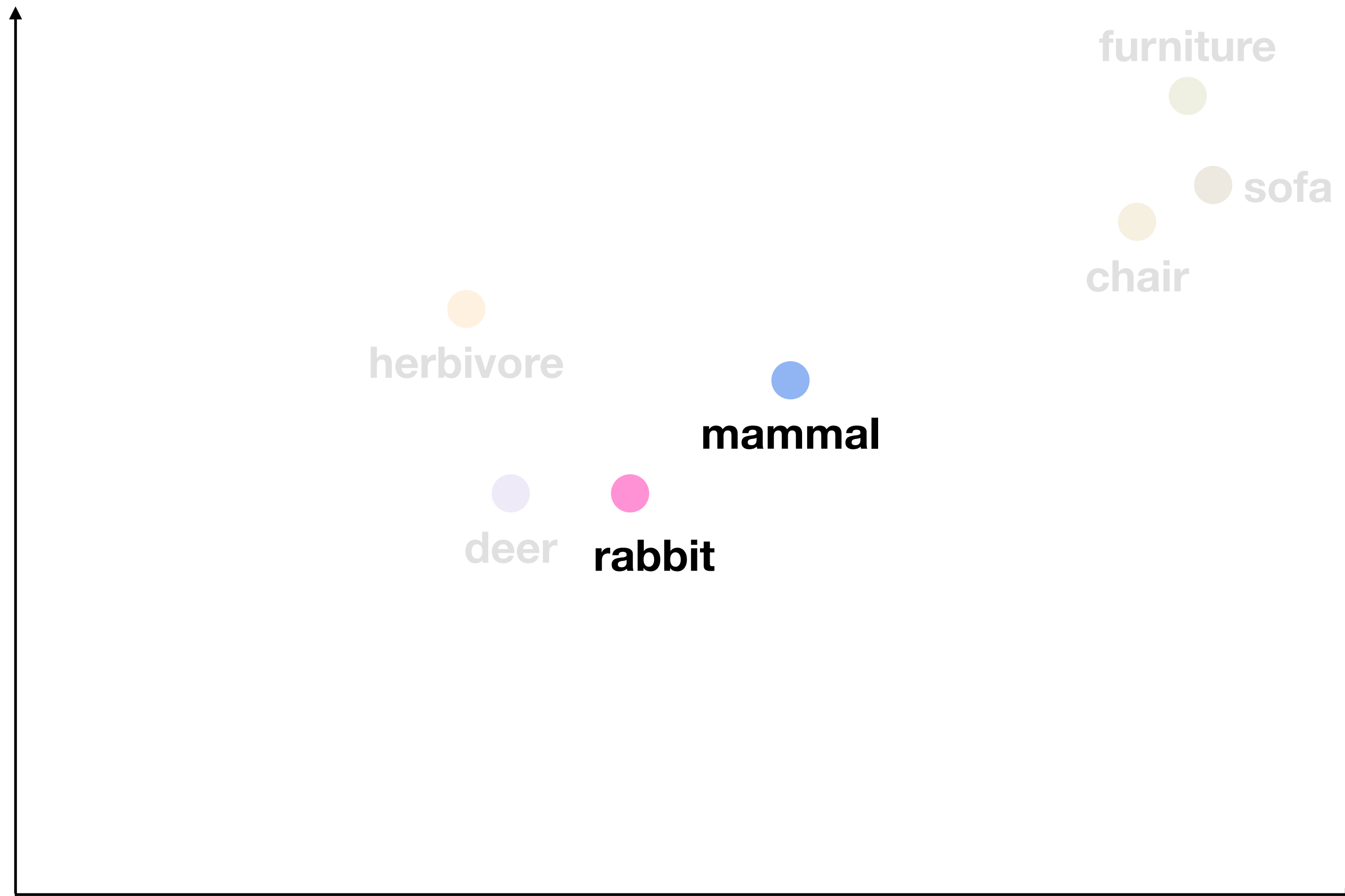
- Commonsense Knowledge.
- Learn the Right Representation.
- Commonsense Knowledge in Pre-trained LMs.
- Benchmark Datasets for Evaluation.

Vector Representation



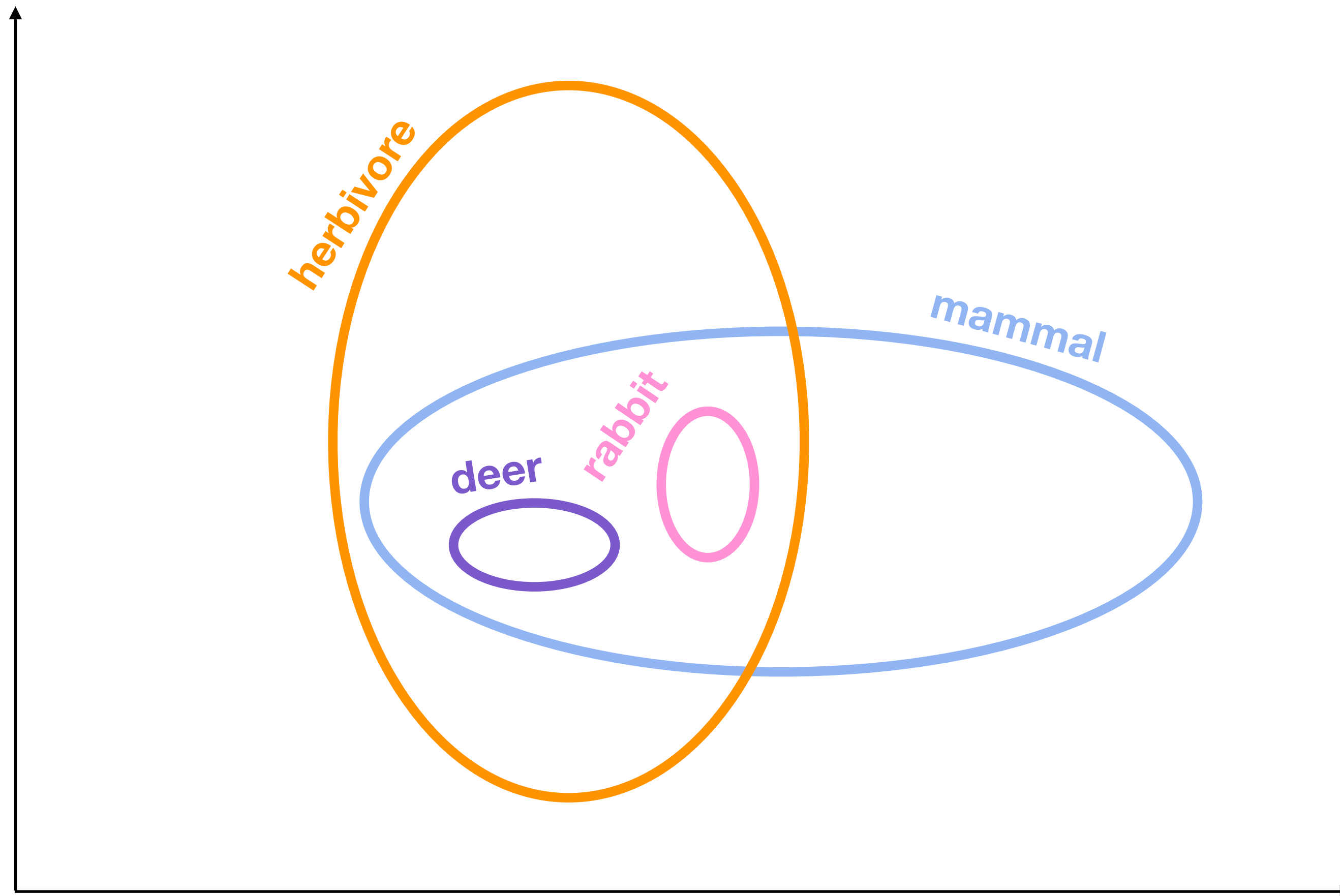
Vector Representation

✗ Region **✗ Asymmetry**



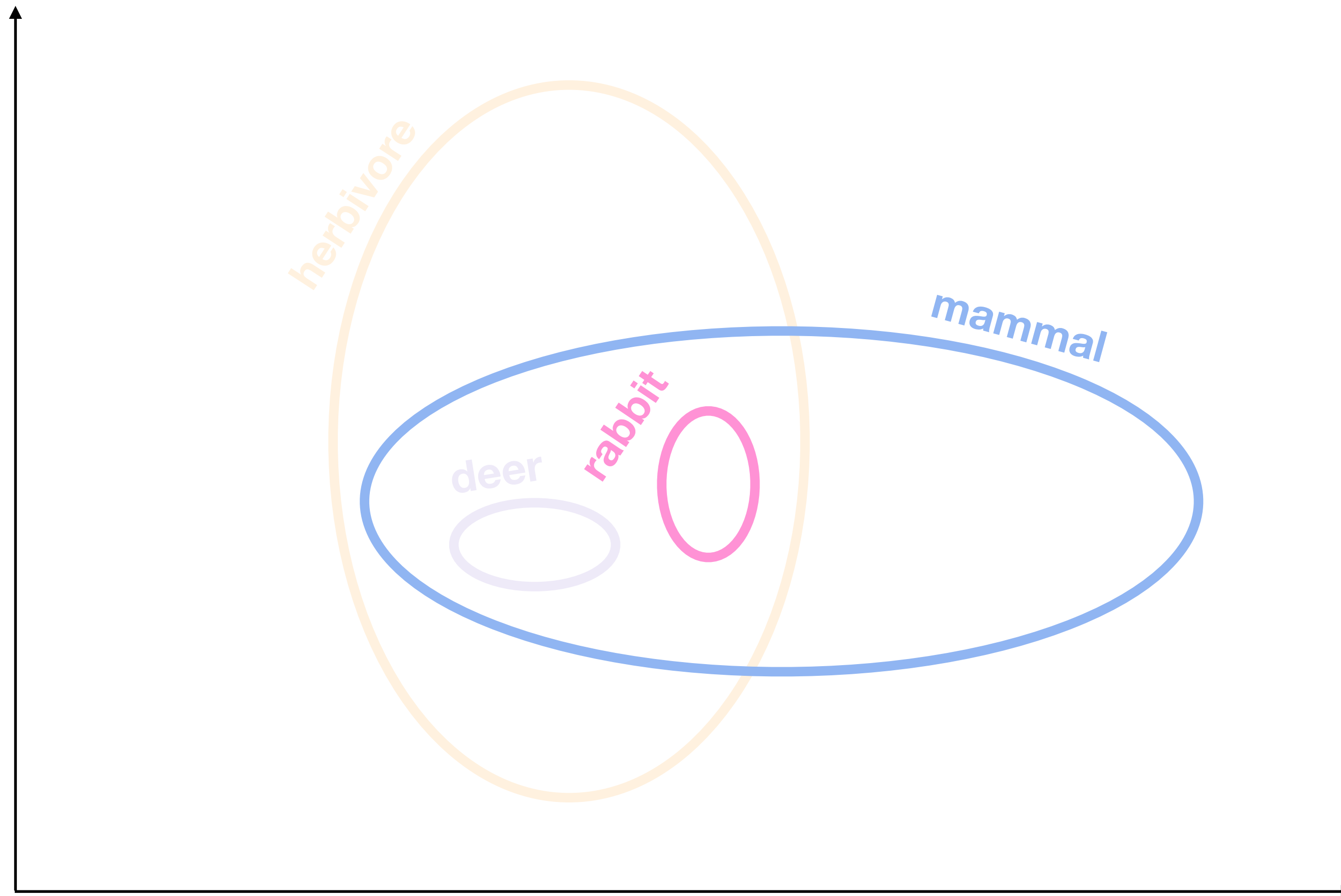
Gaussian Representation

✓ Region ✓ Asymmetry



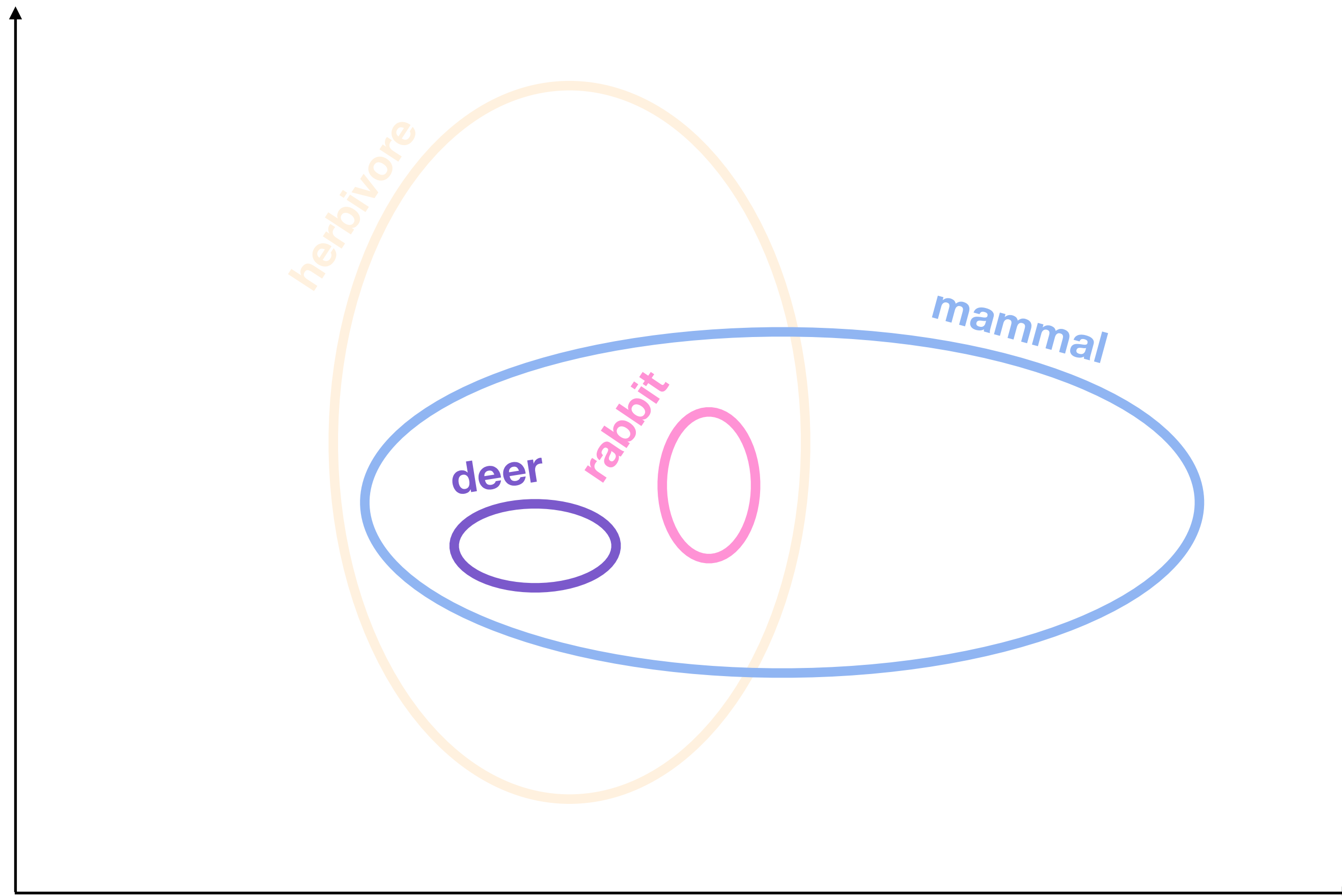
Gaussian Representation

✓ Region ✓ Asymmetry



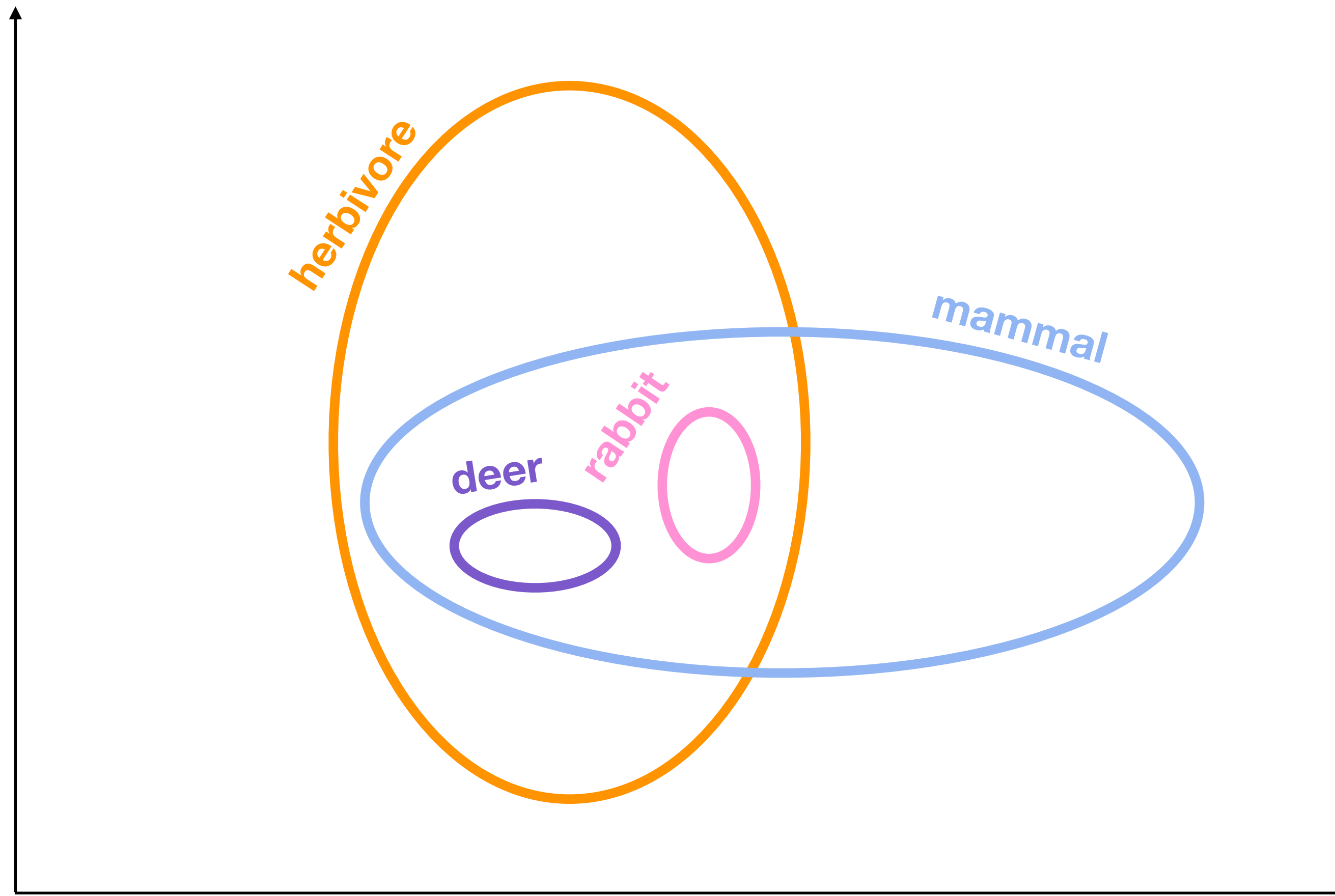
Gaussian Representation

✓ Region ✓ Asymmetry ✓ Disjointness



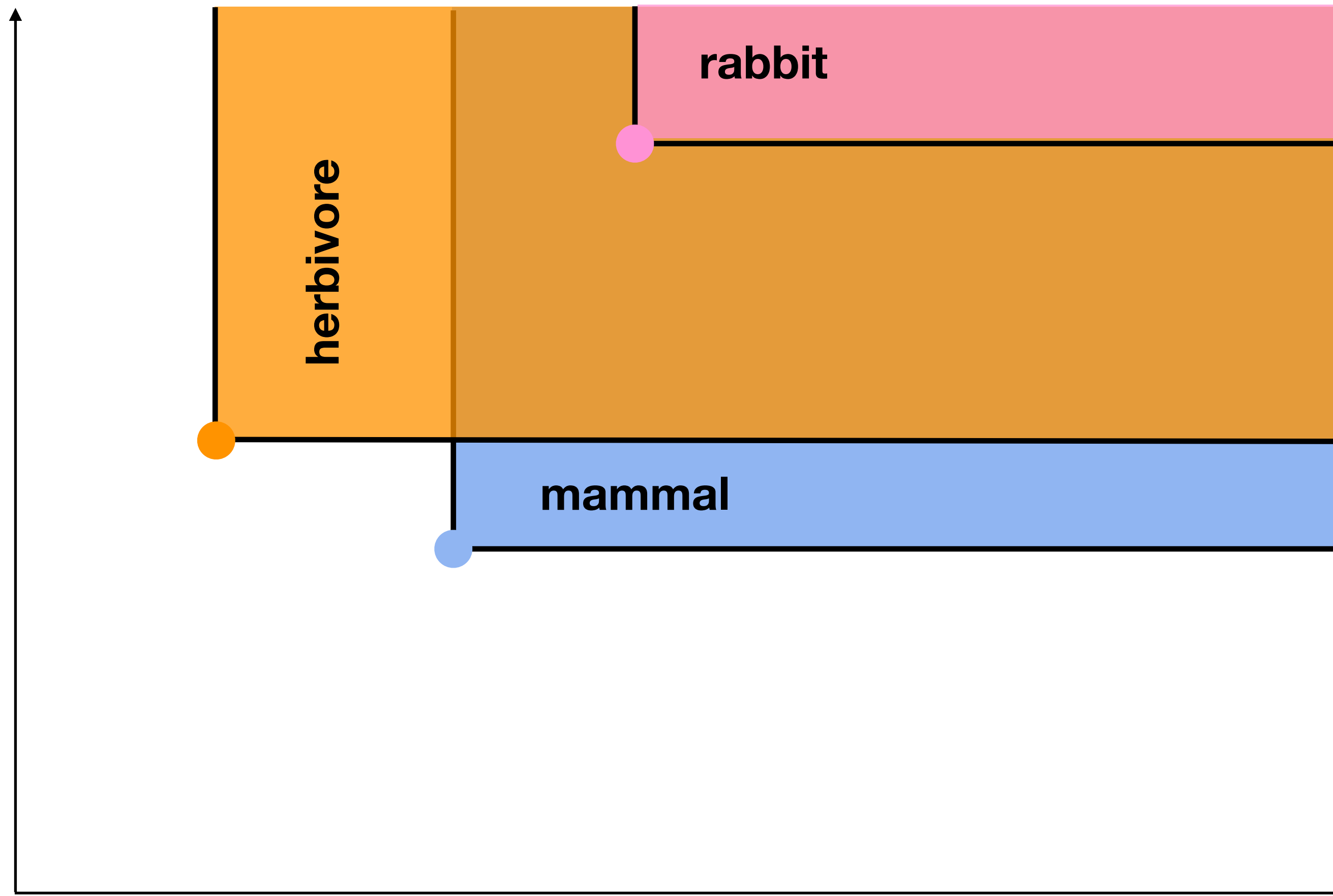
Gaussian Representation

✓ Region ✓ Asymmetry ✓ Disjointness ✗ Closed under intersection



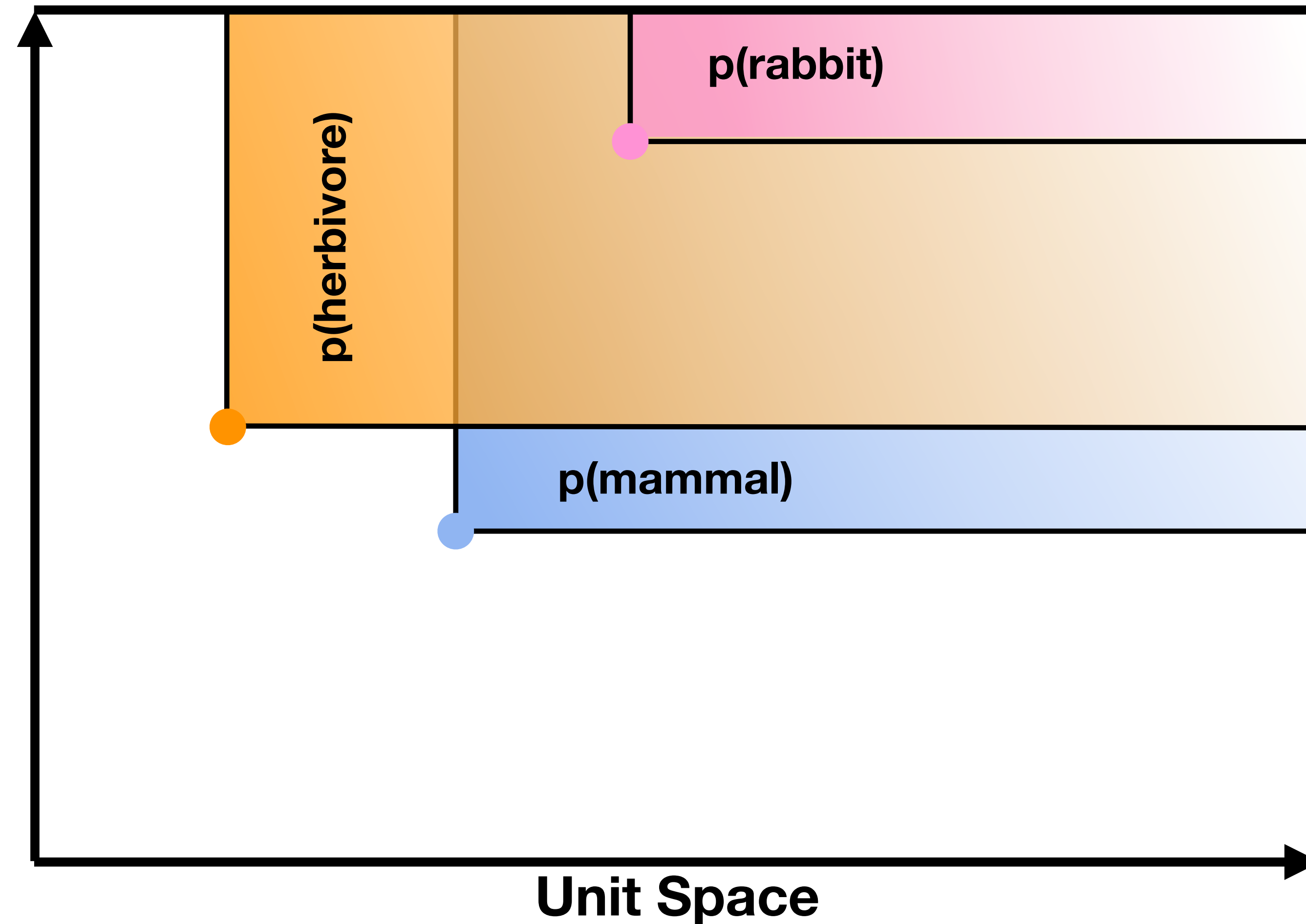
Cone Representation

✓ Region ✓ Asymmetry ✗ Disjointness ✓ Closed under intersection



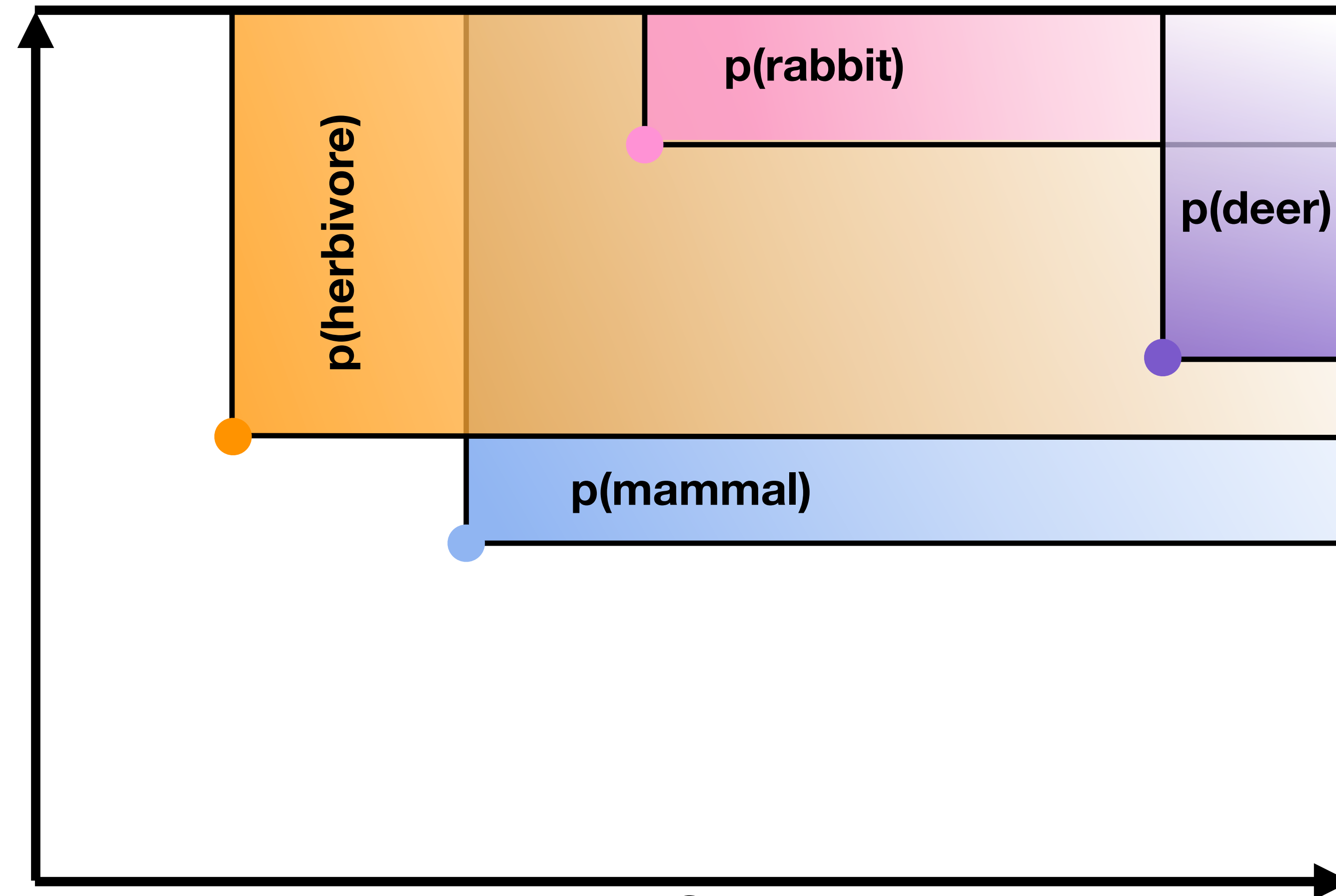
Cone Representation

✓ Region ✓ Asymmetry ✗ Disjointness ✓ Closed under intersection



Cone Representation

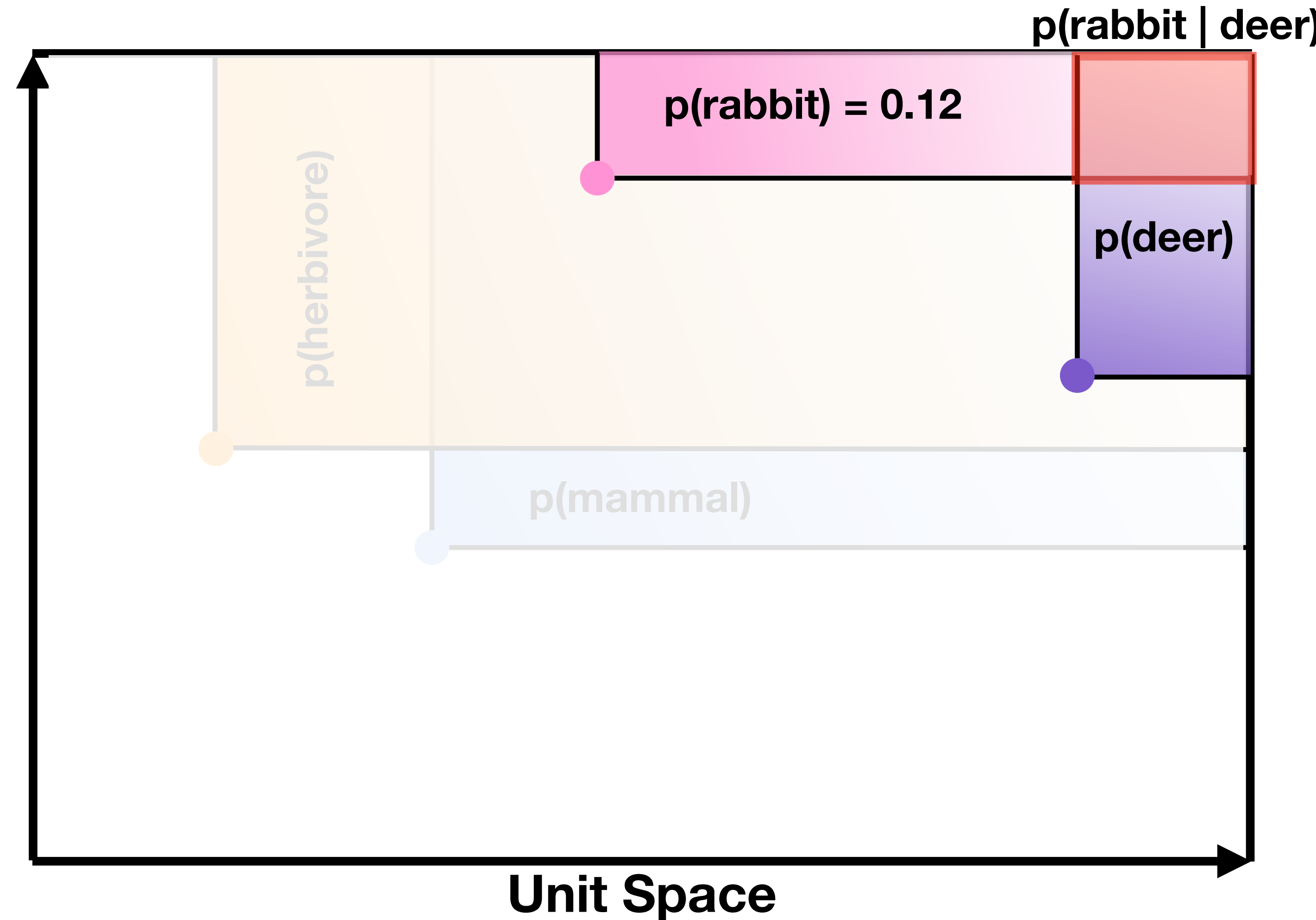
✓ Region ✓ Asymmetry ✗ Disjointness ✓ Closed under intersection



Unit Space

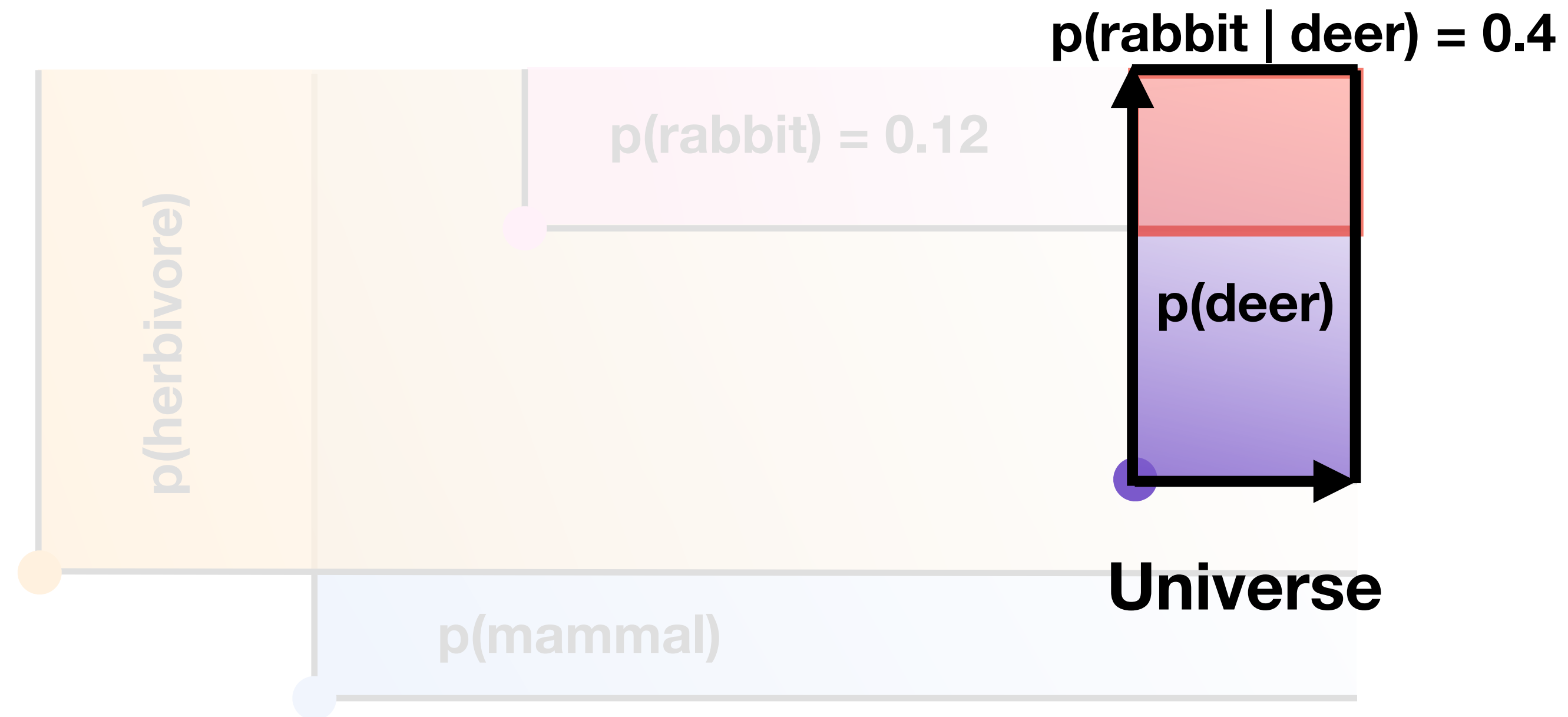
Cone Representation

✓ Region ✓ Asymmetry ✗ Disjointness ✓ Closed under intersection



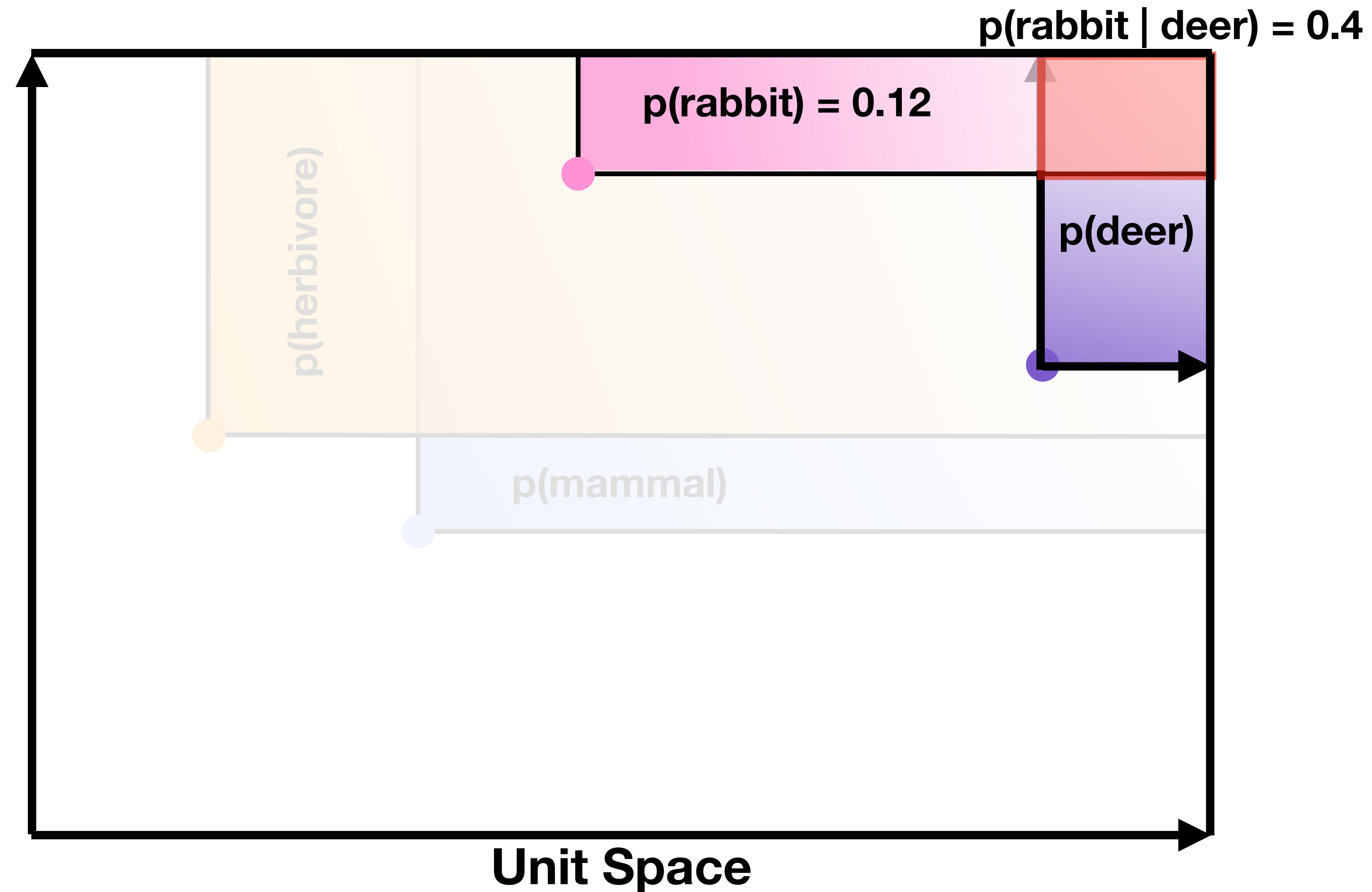
Cone Representation

✓ Region ✓ Asymmetry ✗ Disjointness ✓ Closed under intersection



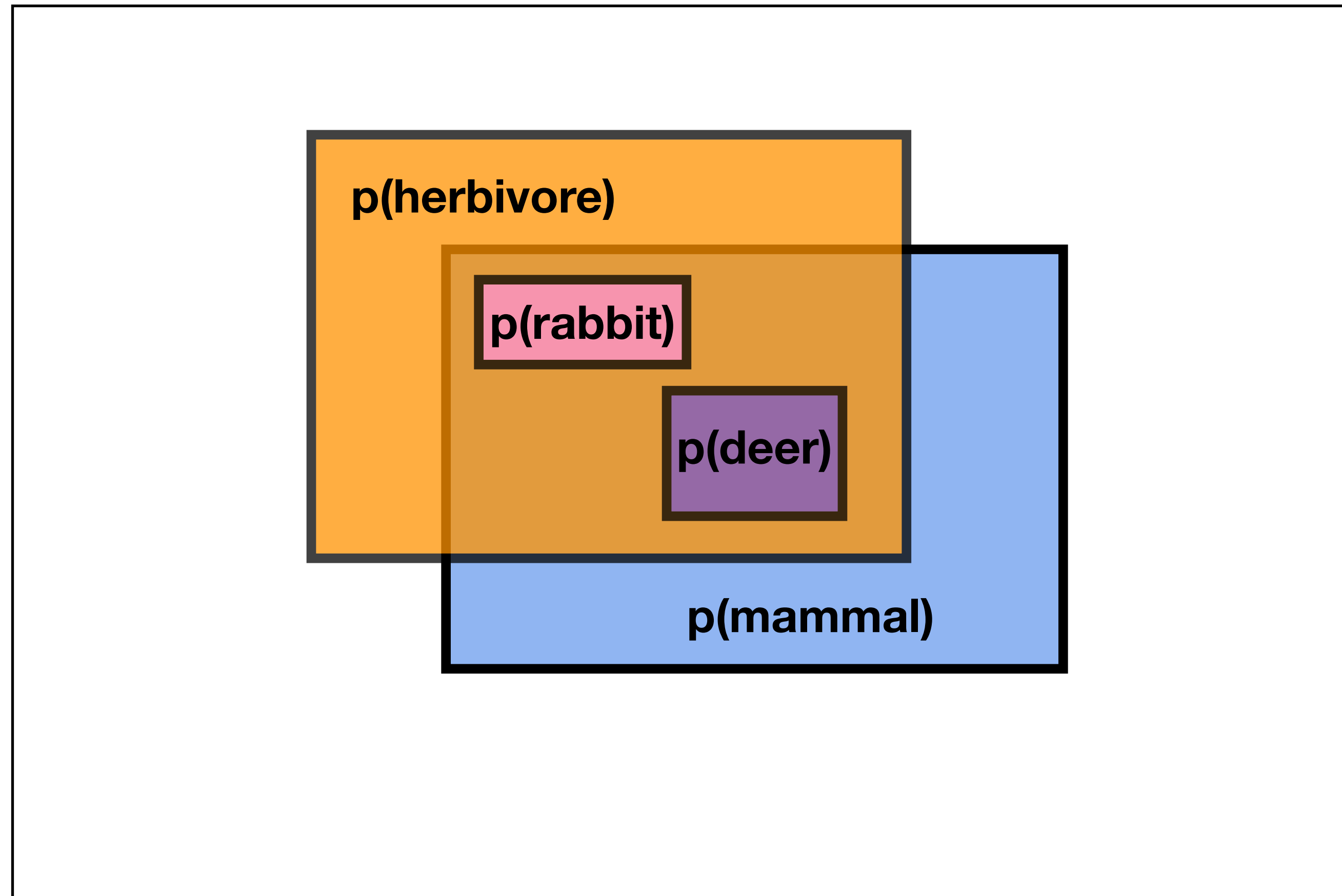
Cone Representation

✓ Region ✓ Asymmetry ✗ Disjointness ✓ Closed under intersection



Box Representation

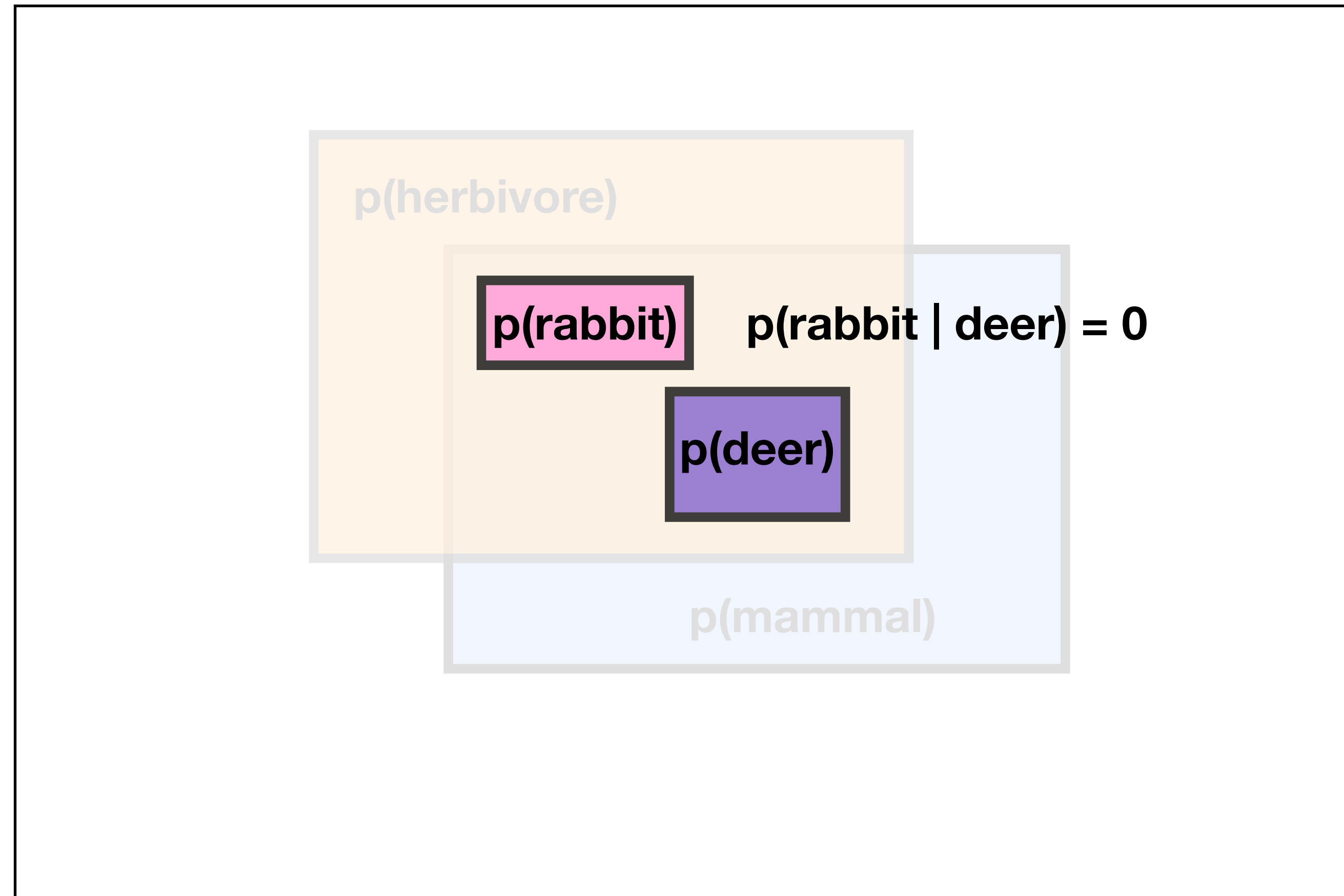
✓ Region ✓ Asymmetry ✓ Disjointness ✓ Closed under intersection



Unit Box

Box Representation

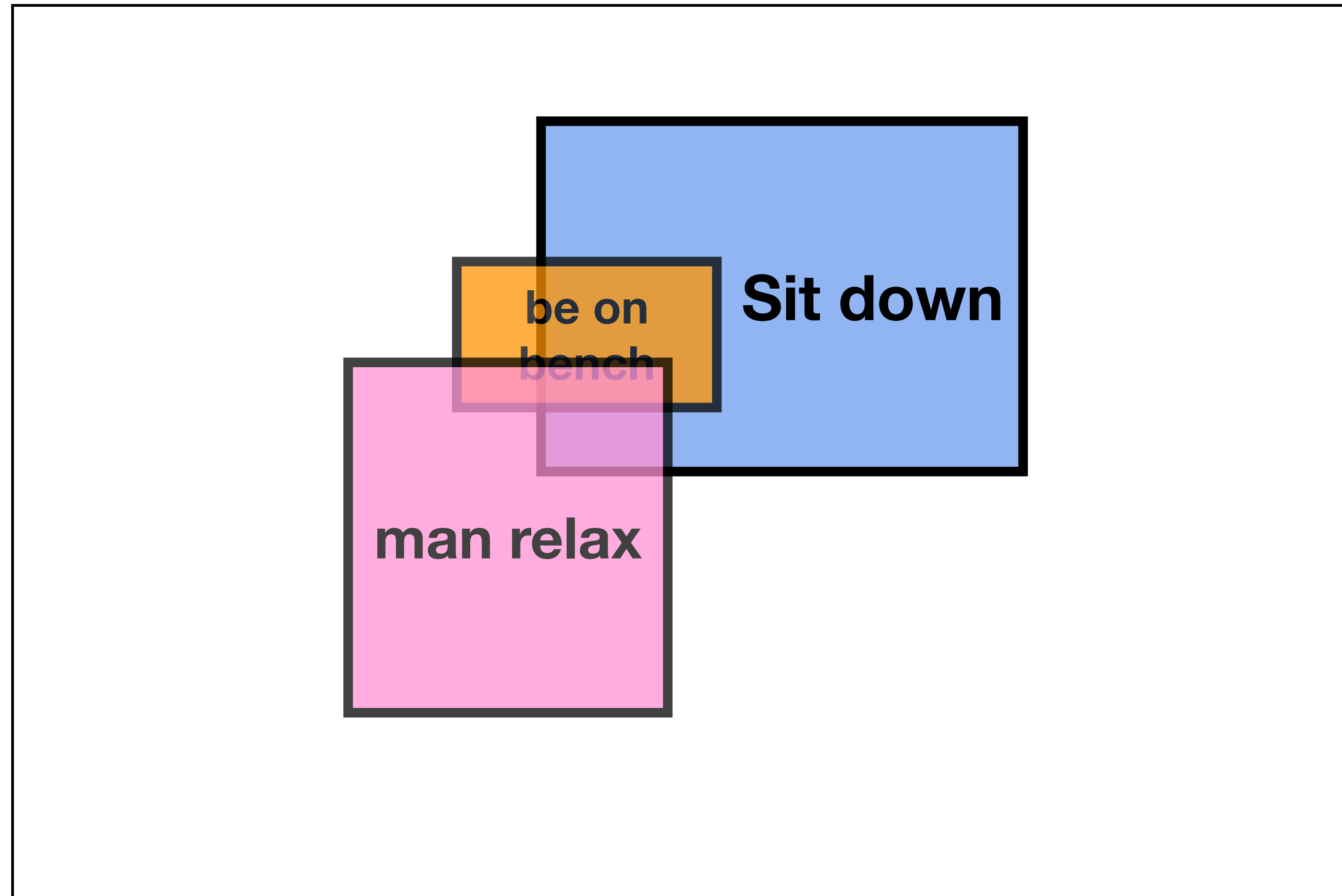
✓ Region ✓ Asymmetry ✓ Disjointness ✓ Closed under intersection



Unit Box

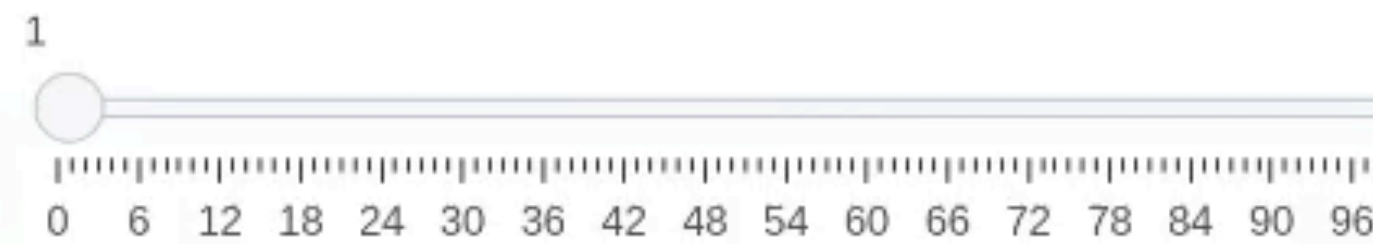
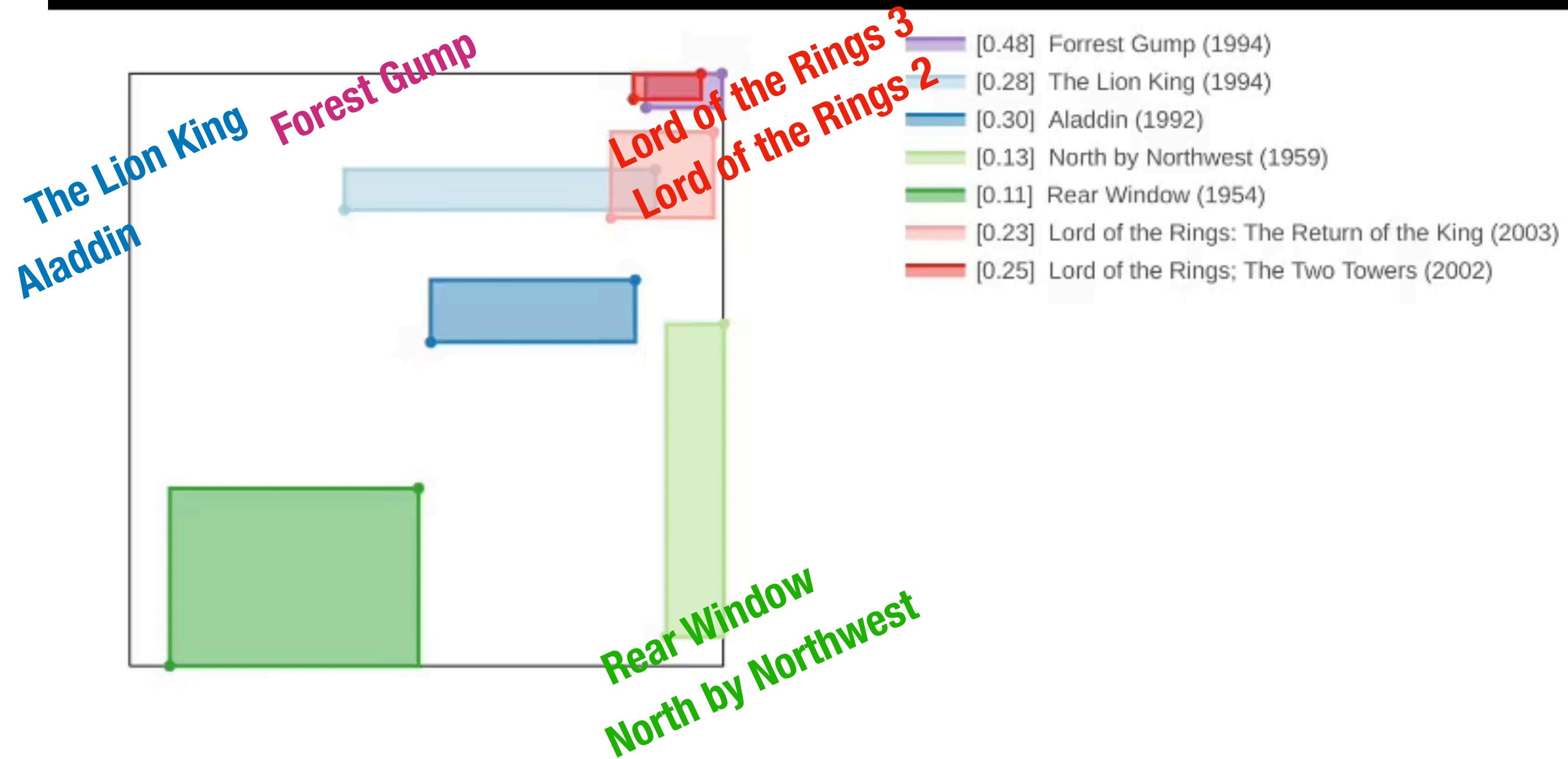
Box Representation

Common Sense

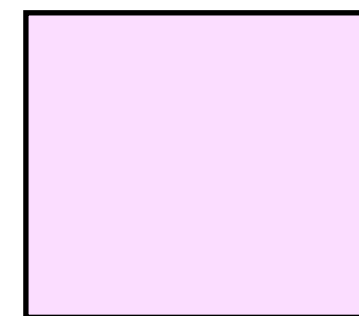
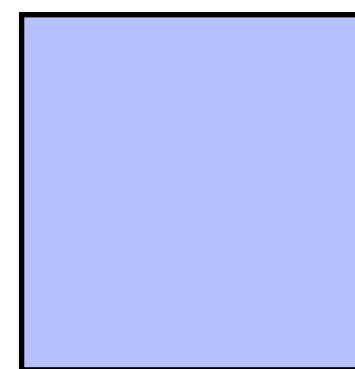
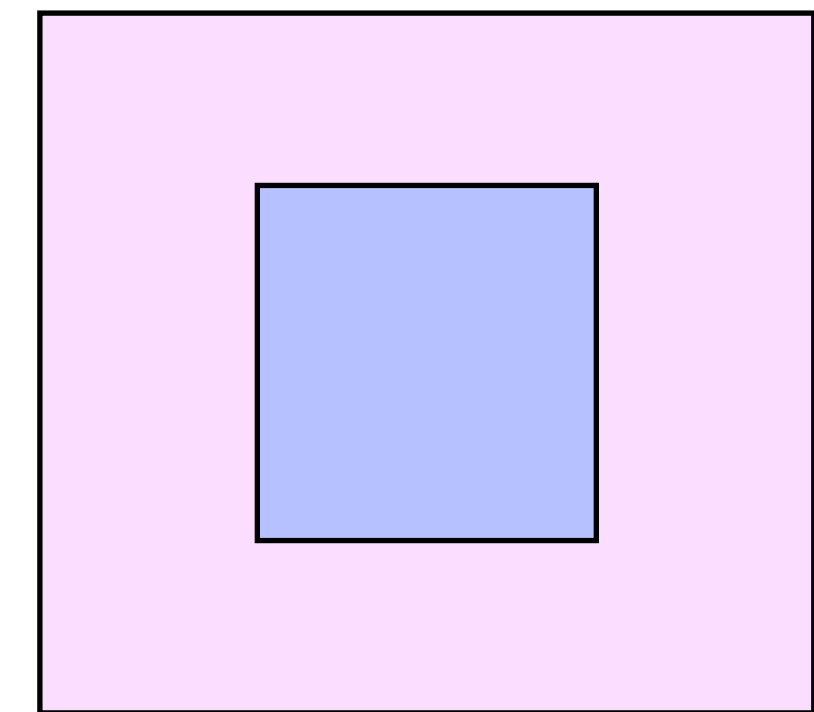
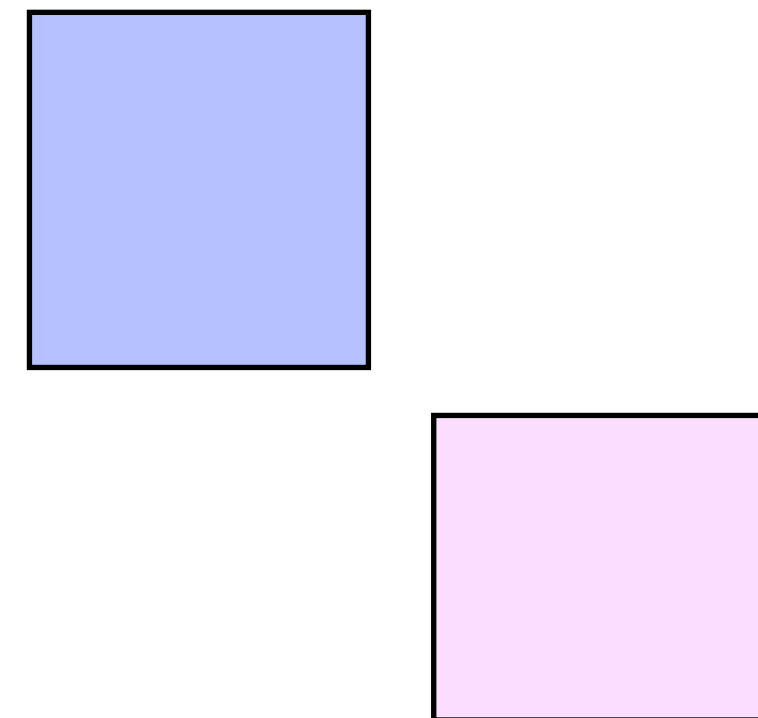
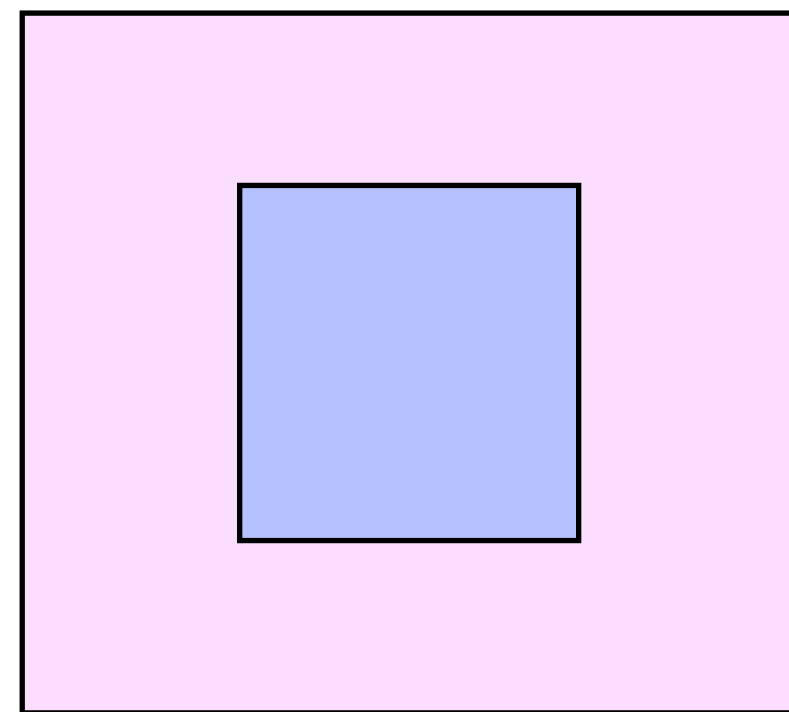


Unit Box

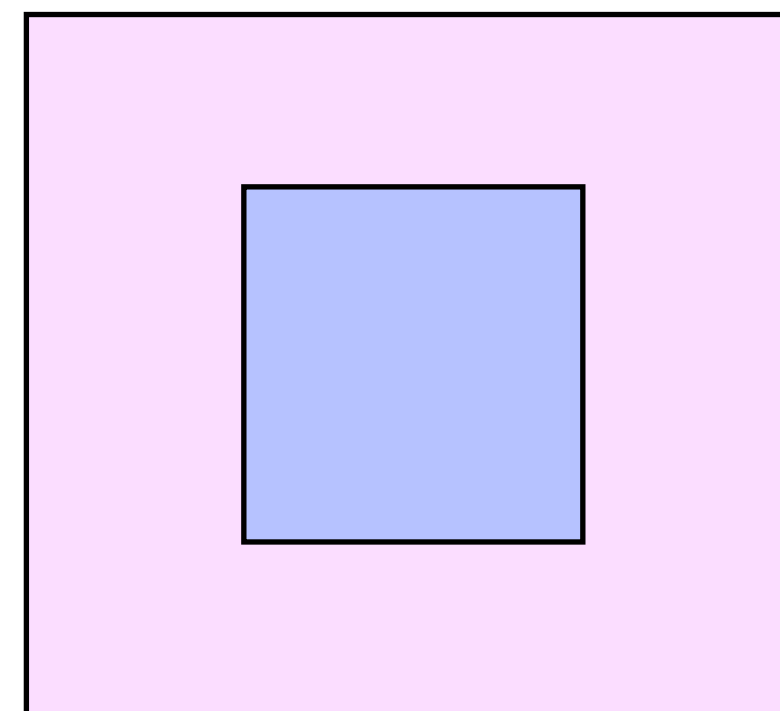
Box Representation Training Video



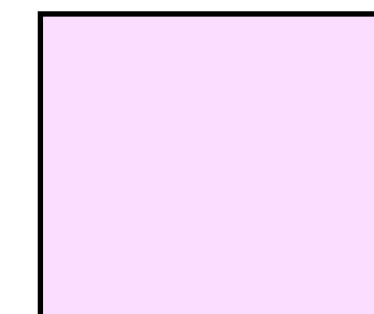
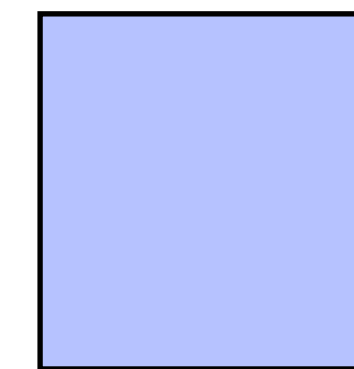
Box Training Difficulty



Current

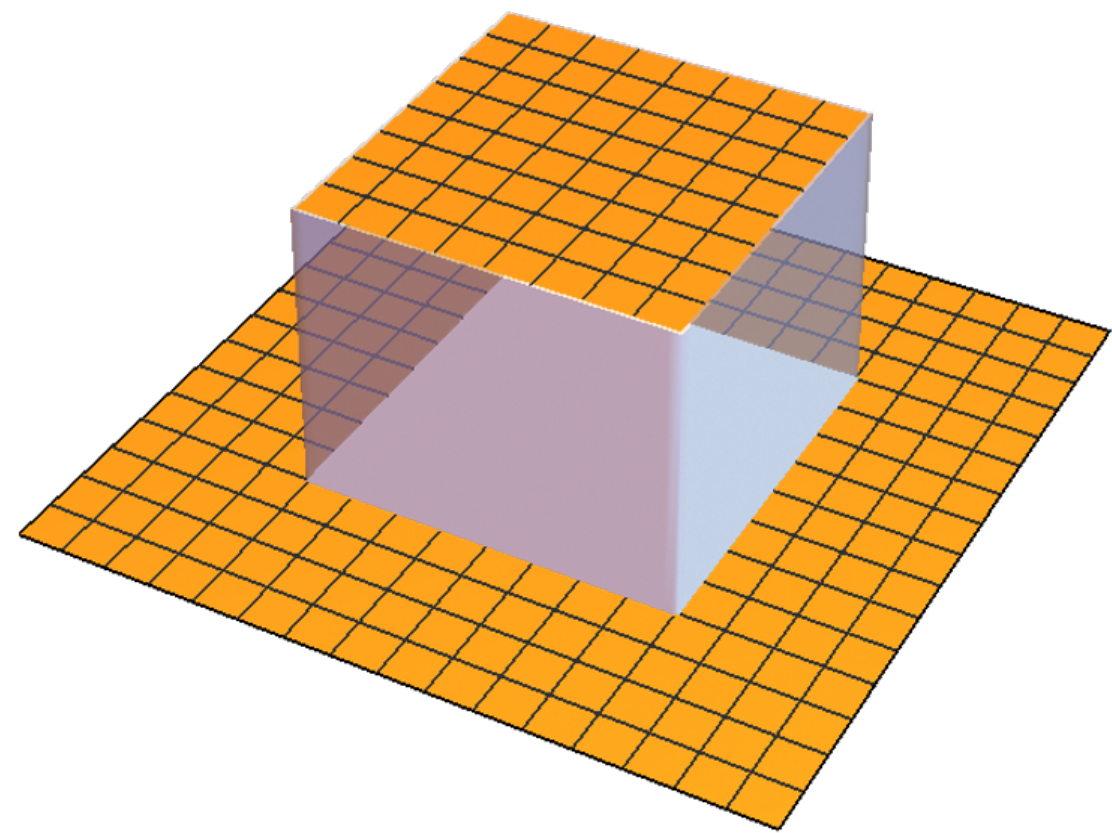


Goal

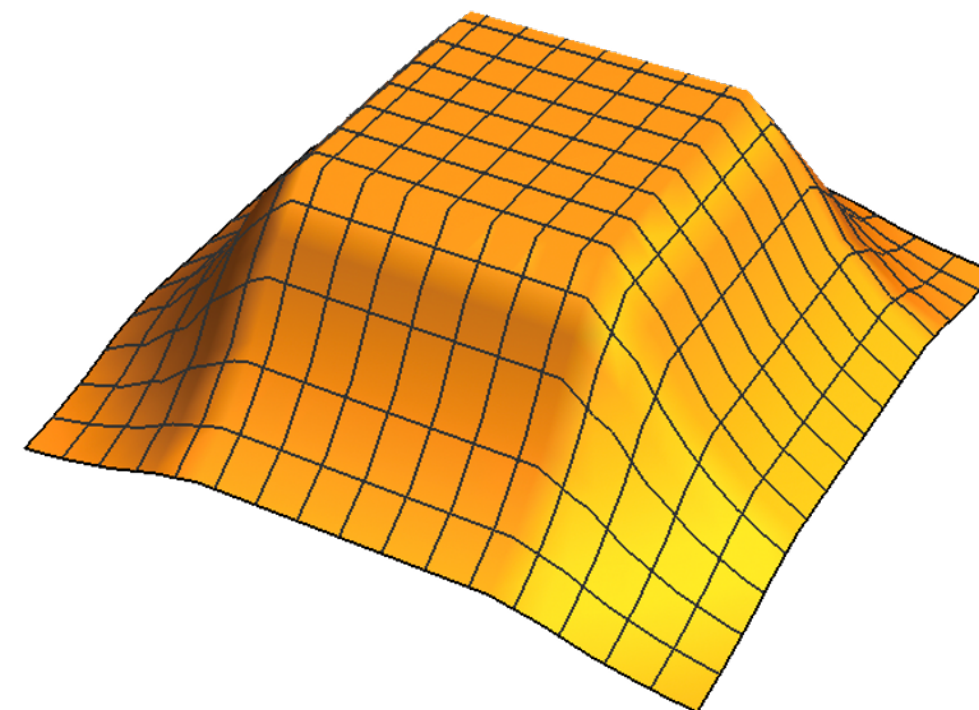


Hard Box Training Result

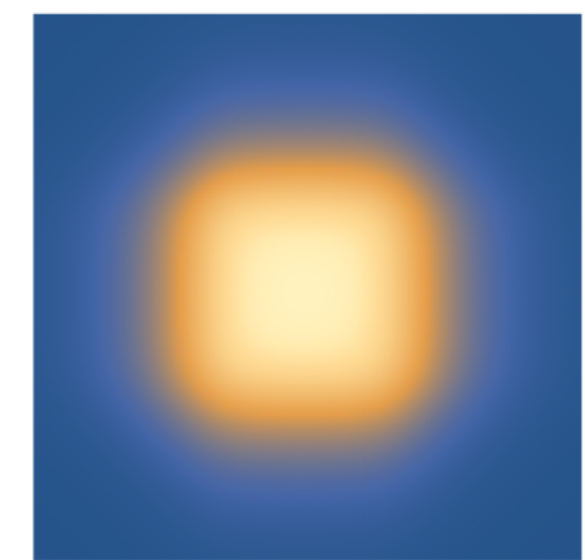
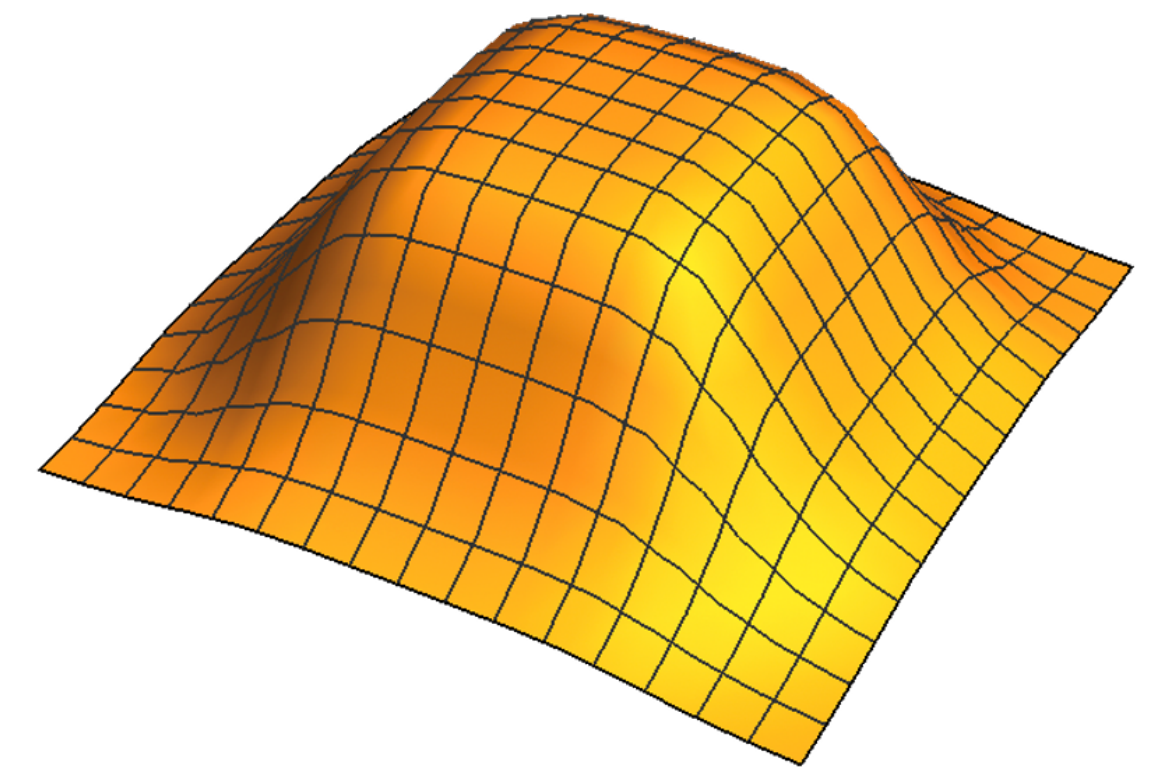
Box Training Loss



Hard Box (ACL 2018)

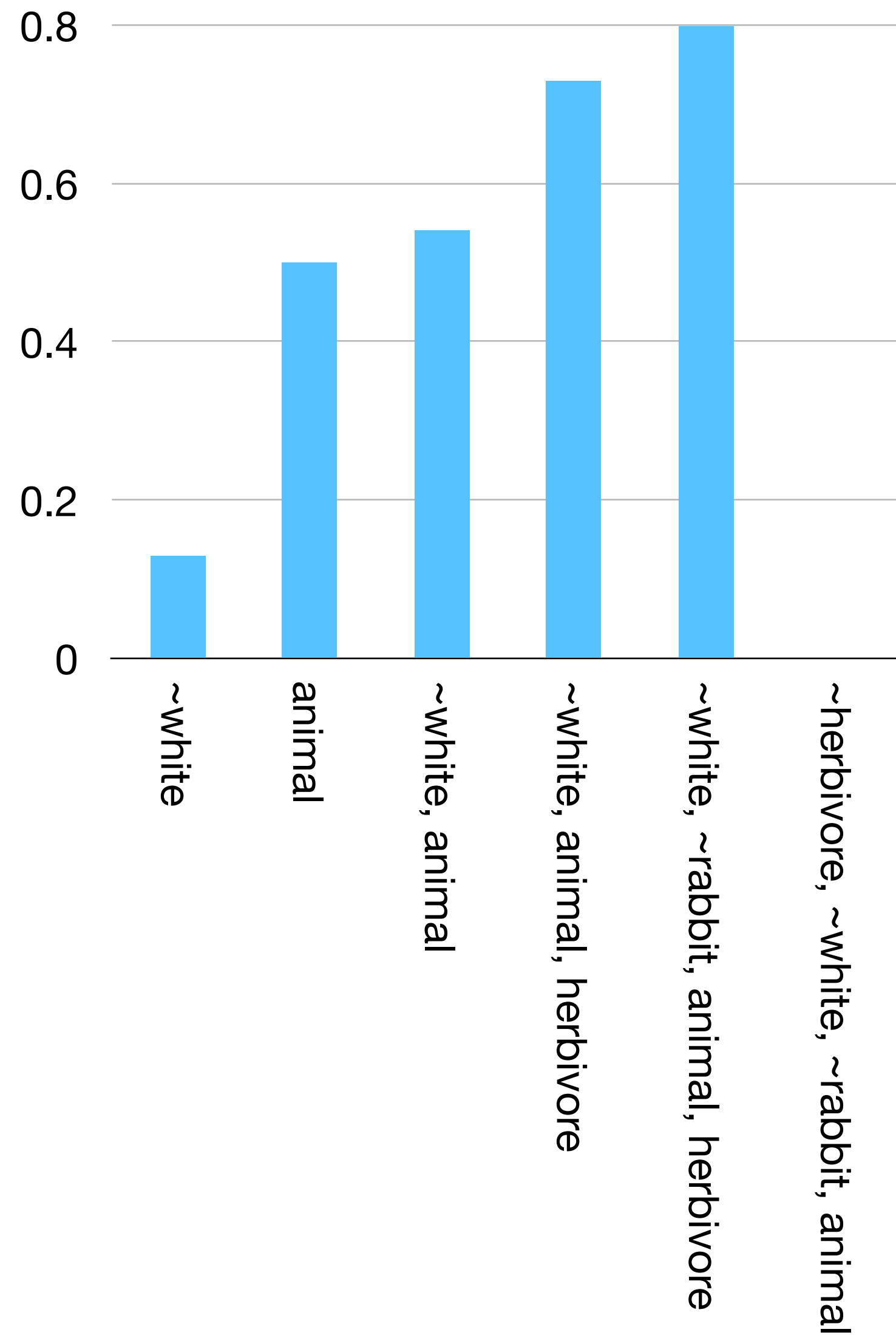


Smoothed Box (ICLR 2019)



Gumbel Box (NeurIPS 2021)

Words examples



	Deer
P(deer)	0.11
~white	0.13
animal	0.50
~white, animal	0.54
~white, animal, herbivore	0.73
~white, ~rabbit, animal, herbivore	0.80
~herbivore, ~white, ~rabbit, animal	0.00

Sentence examples

- Flickr dataset is an entailment dataset containing 45 million image captions.
- Examples

x	p(x)	y	p(y)	p(x y)
person walk	0.11516	blond woman walk down sidewalk	1.6E-04	1.0
person wear clothing	0.43036	adult dance on floor	3.9E-04	0.9
man play percussion instrument	0.00347	drummer	3.4E-03	0.51
man wear jacket	0.03077	snow on ground	5.1E-04	0.31
in basement	4.3E-04	hold instrument	5.9E-03	0.0067

Outline

- Commonsense Knowledge.
- Learn the Right Representation.
- Commonsense Knowledge in Pre-trained LMs.
- Benchmark Datasets for Evaluation.

Do Pre-trained LMs **Already** Capture Commonsense Knowledge?

OpenAI's new language generator GPT-3 is shockingly good—and completely mindless

A robot wrote this entire article. Are you scared yet, human?

We asked GPT-3, OpenAI's powerful new language generator, to write an essay for us from scratch. The assignment? To convince us robots come in peace

- For more about GPT-3 and how this essay was written and edited, please read our editor's note below



New York Times Opinion @nytopinion · Jul 30

"GPT-3 is capable of generating entirely original, coherent and sometimes even factual prose," writes [@fmanjoo](#). "And not just prose — it can write poetry, dialogue, memes, computer code and who knows what else."



Opinion | How Do You Know a Human Wrote This?

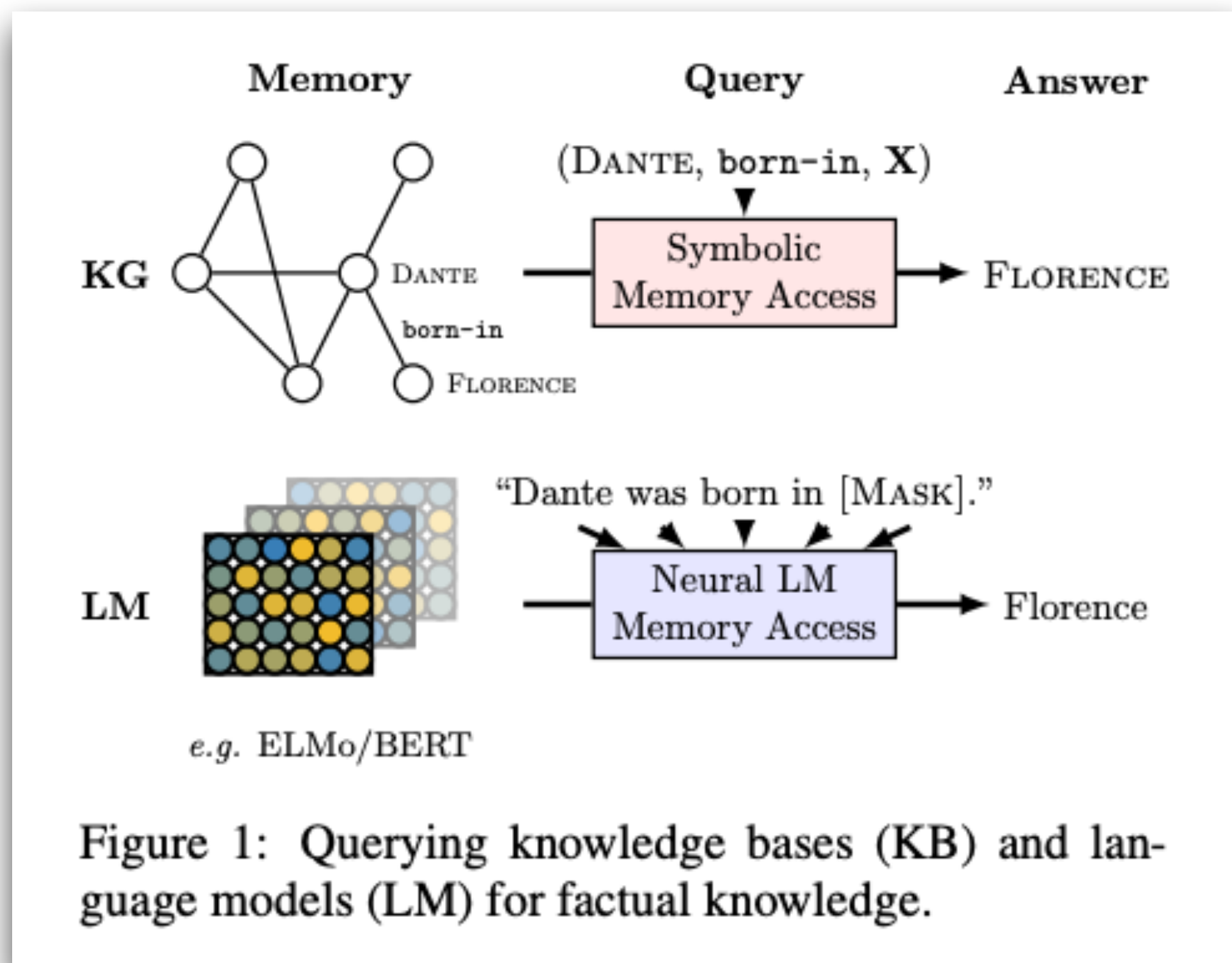
Machines are gaining the ability to write, and they are getting terrifyingly good at it.

[nytimes.com](#)

Do Pre-trained LMs **Already** Capture Commonsense Knowledge?

Commonsense KB relations => Natural language template => Using LMs to query / score

- LAMA: Petroni et al. (EMNLP 2019)
- Feldman et al. (EMNLP 2019)



Candidate Sentence S_i	$\log p(S_i)$
"musician can playing musical instrument"	-5.7
"musician can be play musical instrument"	-4.9
"musician often play musical instrument"	-5.5
"a musician can play a musical instrument"	-2.9

Table 1: Example of generating candidate sentences. Several enumerated sentences for the triple (musician, CapableOf, play musical instrument). The sentence with the highest log-likelihood according to a pretrained language model is selected.

Does the prompt matter?

- Yes! It matters! AutoPrompt (Shin et al., EMNLP 2020)
- Generating gradient guided prompt.

Prompt Type	Original			T-REx		
	MRR	P@10	P@1	MRR	P@10	P@1
LAMA	40.27	59.49	31.10	35.79	54.29	26.38
LPAQA (Top1)	43.57	62.03	34.10	39.86	57.27	31.16
AUTO PROMPT 5 Tokens	53.06	72.17	42.94	54.42	70.80	45.40
AUTO PROMPT 7 Tokens	53.89	73.93	43.34	54.89	72.02	45.57

Do Pre-trained LMs **Already** Capture Commonsense Knowledge?

Properties of Concepts (Weir et al., 2020)

1. Do pre-trained LM correctly **distinguish concepts associated with a given set of properties?**
2. Can pre-trained LMs be used to **list the properties associated with given concepts?**



Do Pre-trained LMs **Already** Capture Commonsense Knowledge?

Properties of Concepts (Weir et al., 2020)

1. Do pre-trained LM correctly **distinguish concepts associated with a given set of properties?**
 - **A ___ has fur.**
 - **A ___ has fur, is big, and has claws.**
 - **A ___ has fur, is big, and has claws, has teeth, is an animal, eats, is brown...**



Do Pre-trained LMs **Already** Capture Commonsense Knowledge?

Properties of Concepts (Weir et al., 2020)

1. Do pre-trained LM correctly **distinguish concepts associated with a given set of properties?**
 - Good performance, RoBERTa > BERT
 - Perceptual (e.g. visual) < non-perceptual (e.g. encyclopaedic or functional).
 - Highly-ranked incorrect answers typically apply to a subset of properties.



Do Pre-trained LMs **Already** Capture Commonsense Knowledge?

Properties of Concepts (Weir et al., 2020)

1. Can pre-trained LMs be used to **list the properties associated with given concepts?**
 - Low correlation with human elicited properties, but coherent and mostly “verifiable by humans”

Context	Human		ROBERTA-L	
	Response	PF	Response	p_{LM}
<i>(Everyone knows that) a bear has ____ .</i>	fur	27	teeth	.36
	claws	15	claws	.18
	teeth	11	eyes	.05
	cubs	7	ears	.03
	paws	7	horns	.02



Can we trust knowledge from LMs?

- **LMs also generate fictitious facts!**

Distributionally-related:

**Barack's Wife Hillary:
Using Knowledge Graphs for Fact-Aware Language Modeling**

Robert L. Logan IV* **Nelson F. Liu†§** **Matthew E. Peters§**
Matt Gardner§ **Sameer Singh***

* University of California, Irvine, CA, USA

† University of Washington, Seattle, WA, USA

§ Allen Institute for Artificial Intelligence, Seattle, WA, USA

{rlogan, sameer}@uci.edu, {mattg, matthewp}@allenai.org, nliu@cs.washington.edu

Syntactically-similar:

**Negated and Misprimed Probes for Pretrained Language Models:
Birds Can Talk, But Cannot Fly**

Nora Kassner, Hinrich Schütze

Center for Information and Language Processing (CIS)

LMU Munich, Germany

kassner@cis.lmu.de








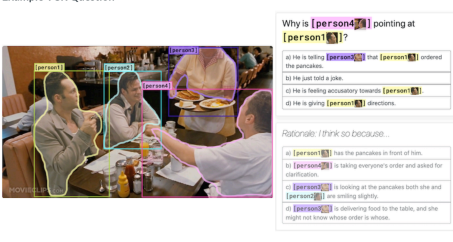
Outline

- Commonsense Knowledge.
- Learn the Right Representation.
- Commonsense Knowledge in Pre-trained LMs.
- Benchmark Datasets for Evaluation.








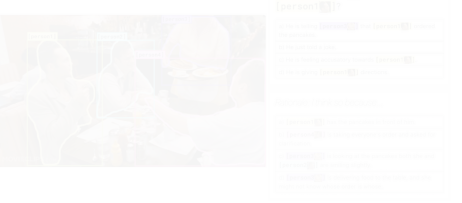
Benchmark Evaluation Dataset

- **Besides probing the model using commonsense knowledge bases, are there standard commonsense benchmark datasets to evaluate the model?**
 - **Question Answering.**
 - **Natural Language Inference.**
 - **Coreference Resolution.**
 - **...**


Benchmark Evaluation Dataset

	Task Type	Domain	Example	Gap to Human Performance
	Multi Choice QA	Grounded commonsense.		95.6% - 93.85% = 1.75%
	Multi Choice Selection	Abductive Reasoning	<p>Obs1: Jenny was addicted to sending text messages.</p> <p>Obs2: Jenny narrowly avoided a car accident.</p> <p>Hyp1: Since her friend's texting and driving car accident, Jenny keeps her phone off while driving.</p> <p>Hyp2: Jenny was looking at her phone while driving so she wasn't paying attention.</p>	92.90% - 89.70% = 3.2%
	Multi Choice QA	Reading Comprehension	<p>Example 1.</p> <p>Paragraph: It's a very humbling experience when you need someone to dress you every morning, for your shoes, and put your hair up. Every mental task takes an unprecedented amount of effort. It made me appreciate Dan even more. But anyway I shan't dwell on this (I'm not dying after all) and not let it detract from my lovely 5 days with my friends visiting from Jersey.</p> <p>Question: What's a possible reason the writer needed someone to dress him every morning?</p> <p>Option1: The writer doesn't like putting effort into these tasks.</p> <p>Option2: The writer has a physical disability.</p> <p>Option3: The writer is bad at doing his own hair.</p>	94% - 91.79% = 2.3%
	<p>Example HellaSwag Question</p> <p>A woman is outside with a bucket and a dog. The dog is running around trying to avoid a bath. She</p> <ul style="list-style-type: none"> a) rinses the bucket off with soap and blow dries the dog's head. b) uses a hose to keep it from getting soapy. c) gets the dog wet, then it runs away again. d) gets into the bath tub with the dog. 			= 4.77%
				= 4.95%
				2.72%
		Multi Choice QA	Vision & Language	

Benchmark Evaluation Dataset

	Task Type	Domain	Example	Gap to Human Performance
	Multi Choice QA	Grounded commonsense.	Example HellaSwag Question A woman is outside with a bucket and a dog. The dog is running around trying to avoid a bath. She <ul style="list-style-type: none"> • it throws the bucket off with soap and then dries the dog's head. • it uses a hose to keep it from getting soapy. • it gets the dog wet, then it runs away again. • it gets into the bath tub with the dog. 	95.6% - 93.85% = 1.75%
	Multi Choice Selection	Abductive Reasoning	Q1:1. Jerry was addicted to reading bad news. Q1:2. Jerry recently avoided a car accident. Q1:3. Steve lost her friend's wedding and driving car accident. Jerry keeps her phone off while driving. Q1:4. Jerry was looking at her phone while driving so she wasn't paying attention.	92.90% - 89.70% = 3.2%
	Multi Choice QA	Reading Comprehension	Example 1 The first paragraph of the passage states that the author is a scientist who has spent most of his life studying the universe. He is particularly interested in the study of black holes and has spent many years of his life trying to understand them. He has written several books on the subject and has given many lectures at universities around the world. He is currently working on a new project that he believes will revolutionize our understanding of the universe.	94% - 91.79% = 2.3%
	Multi Choice QA	Naive physical reasoning	Example You need to break a window. Which object would you rather use? <ul style="list-style-type: none"> • a metal stool • a brick • a bottle of water 	94.9% - 90.13% = 4.77%
	Multi Choice QA	Social commonsense	Example Social IQa Question In the urban park, Peter placed a chair on the bench to the bench with the large other. How would other the be react? <ul style="list-style-type: none"> • getting the chair • it would not be able to support • it would not be able to sit 	88.1% - 83.15% = 4.95%
	Multi Choice Selection	Coreference resolution	Reference: Barbara had the greatest means to afford a new car with money she had saved. She had a high paying job. Context: Barbara Context: Barbara	94% - 91.28% = 2.72%
	Multi Choice QA	Vision & Language		85% - 77.79% = 2.3%



Generative Evaluation 

Benchmark Evaluation Dataset

Generative Evaluation

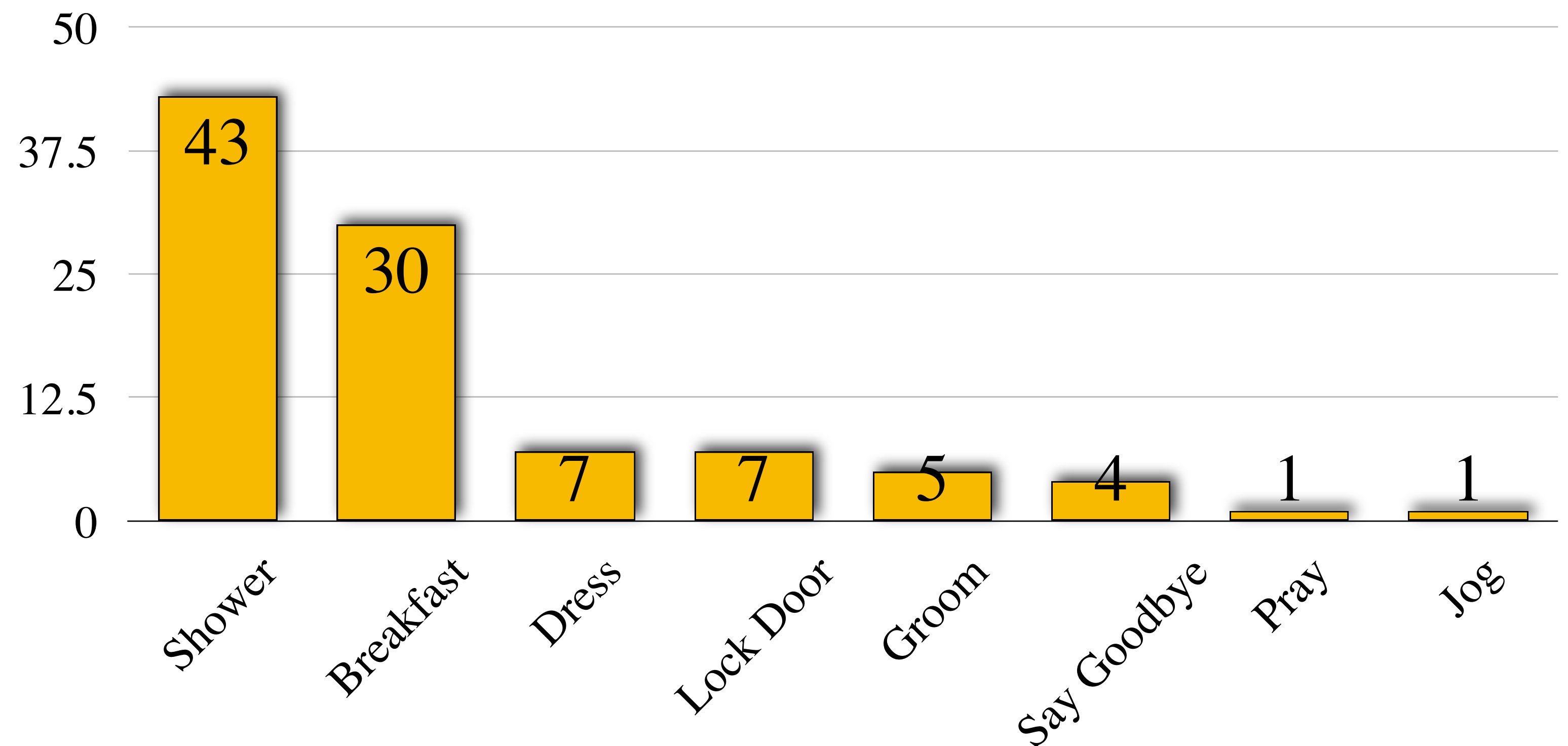


◆ **ProtoQA (EMNLP 2020)**: dataset that captures prototypical situation.

✓ Multiple correct answers.

✓ Scores for each answer.

Name something that people usually do before they leave the house for work?

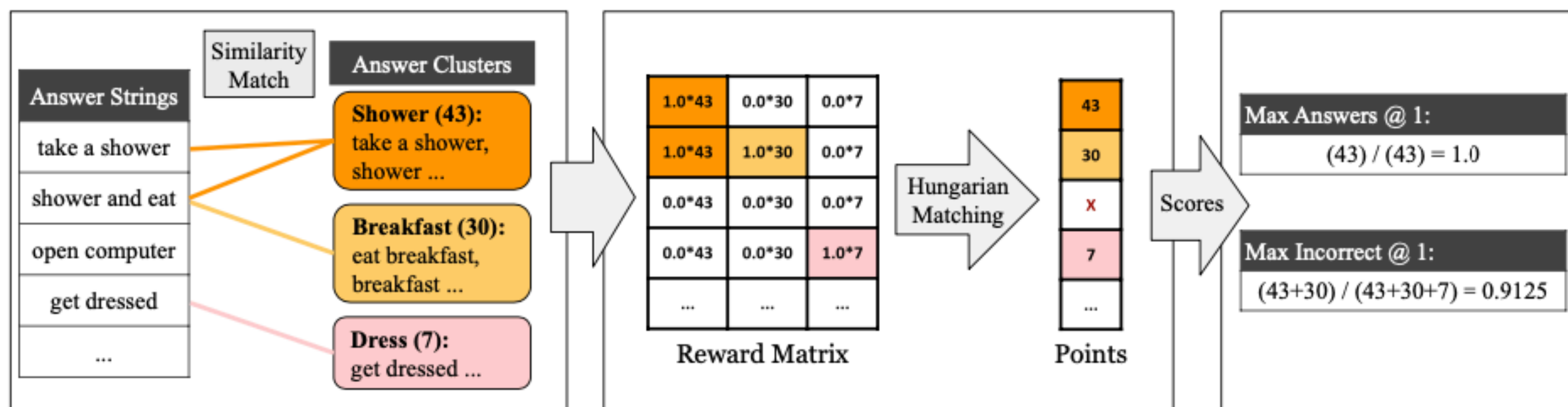


ProtoQA (EMNLP 2020)

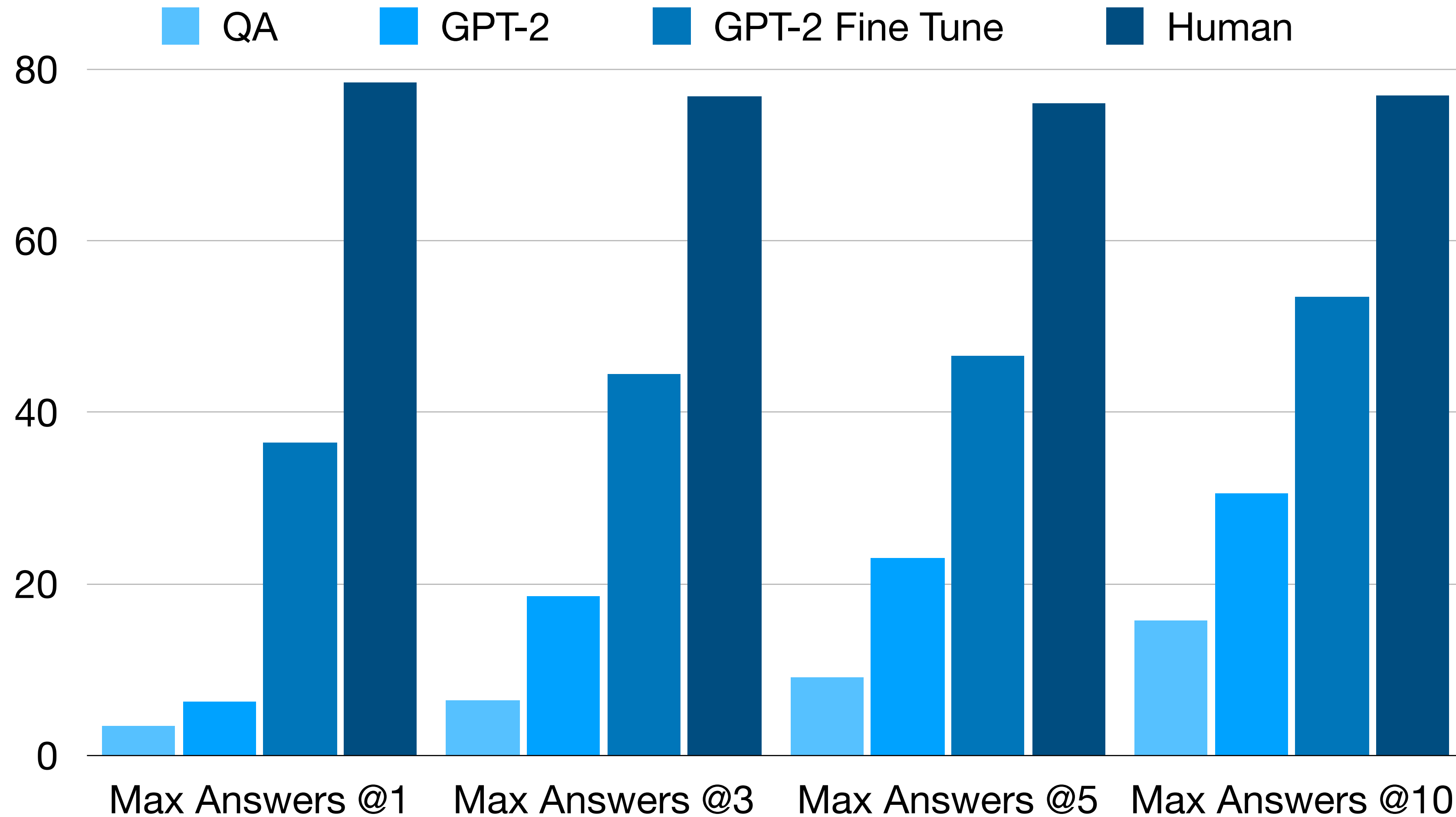
◆ Generative Evaluation

- ✓ Evaluate multiple correct answers generative by the model.
- ✓ Reward models with correct ranking of answer list.
- ✓ Reward models with higher coverage of answer list.

Name something that people usually do before they leave for work.



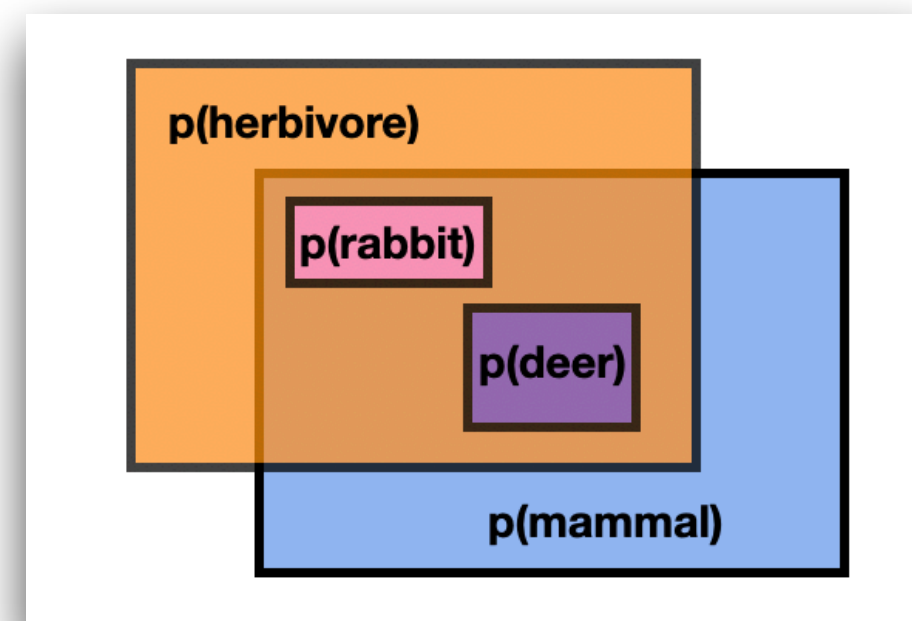
Results



Numbers reported are percentage of perfect score, i.e. answering with a list with an element from each answer cluster in decreasing order would yield 100.

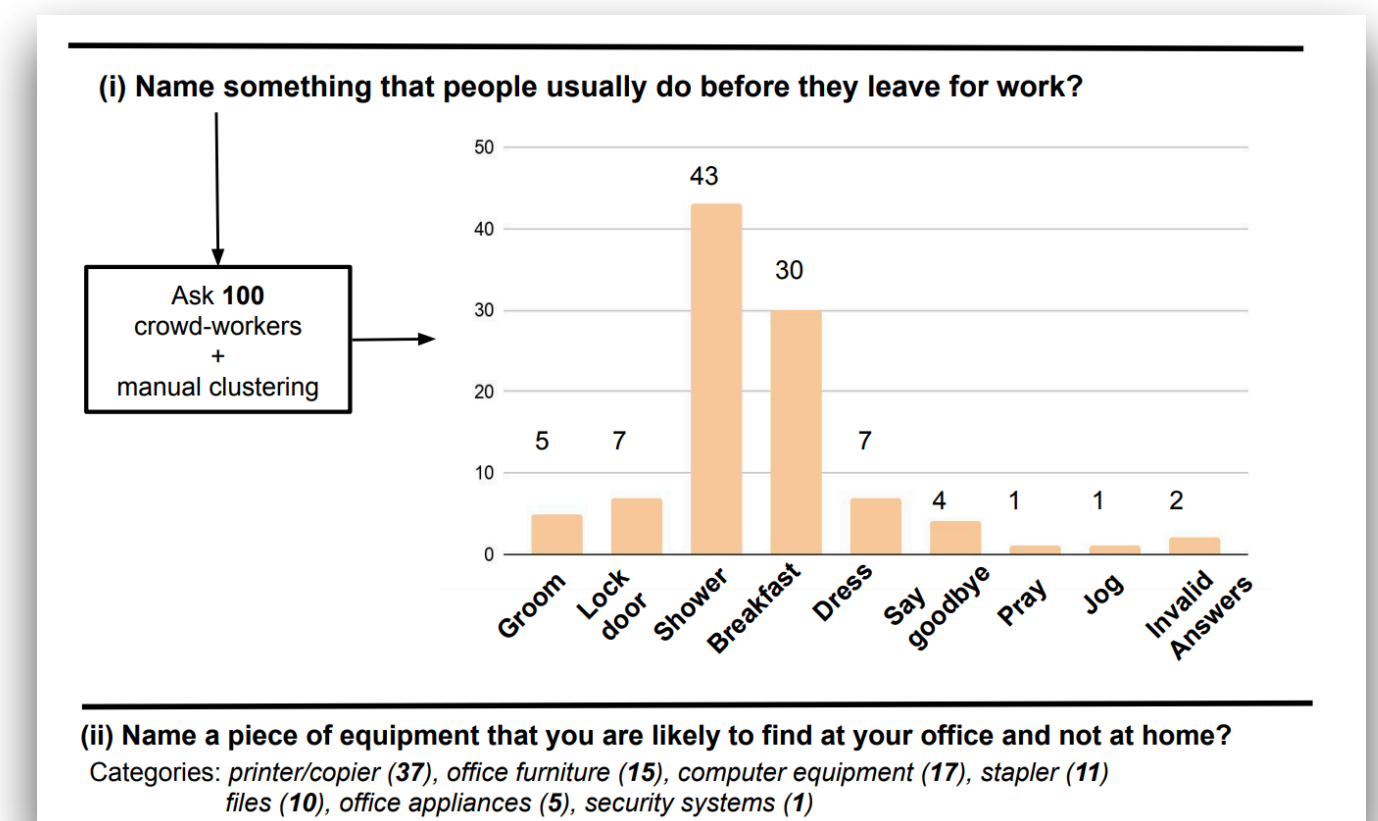
Summary

- Commonsense Knowledge.
- Learn the Right Representation.
- Commonsense Knowledge in Pre-trained LMs.
- Benchmark Datasets for Evaluation.



Candidate Sentence S_i	$\log p(S_i)$
"musician can playing musical instrument"	-5.7
"musician can be play musical instrument"	-4.9
"musician often play musical instrument"	-5.5
"a musician can play a musical instrument"	-2.9

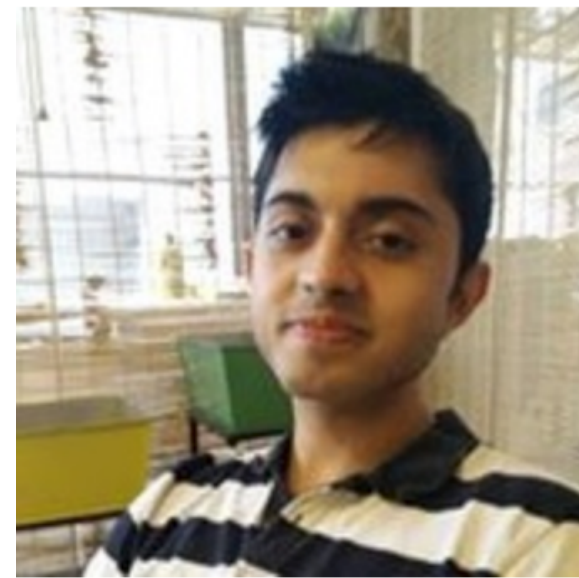
Table 1: Example of generating candidate sentences. Several enumerated sentences for the triple (musician, CapableOf, play musical instrument). The sentence with the highest log-likelihood according to a pretrained language model is selected.



Thanks to all the collaborators!



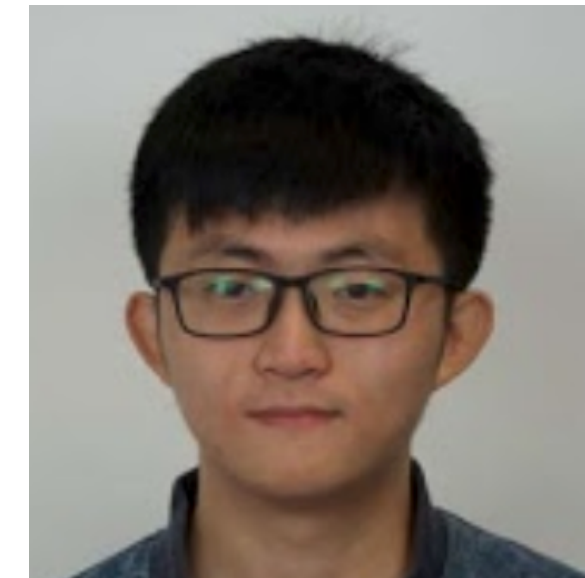
Dan Le



Shikhar Murty



Shib Sankar Dasgupta



Dongxu Zhang



Rajarshi Das



Luke Vilnis



Michael Boratko



Tim O'Gorman



Andrew McCallum