



Fairness, Accountability, and Uncertainty

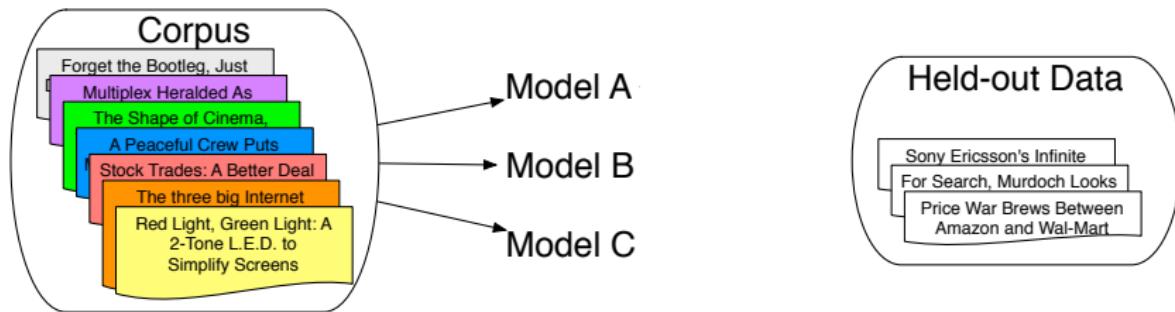
Jordan Boyd-Graber
University of Maryland
INTERPRETABILITY

Slides adapted from Juan Ribeiro

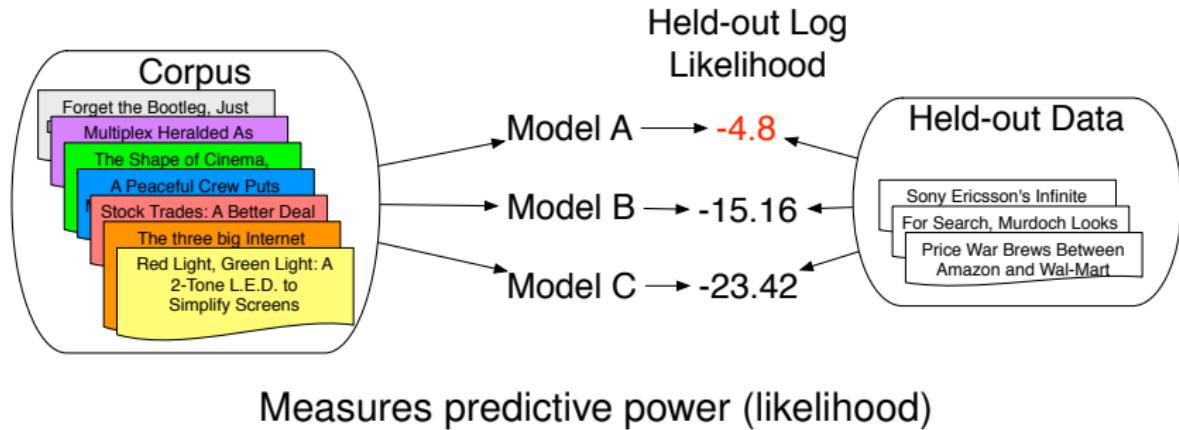
Interpretability

- We believe that interpretability is important
- But need to be able to measure
- Differences for supervised and unsupervised ML

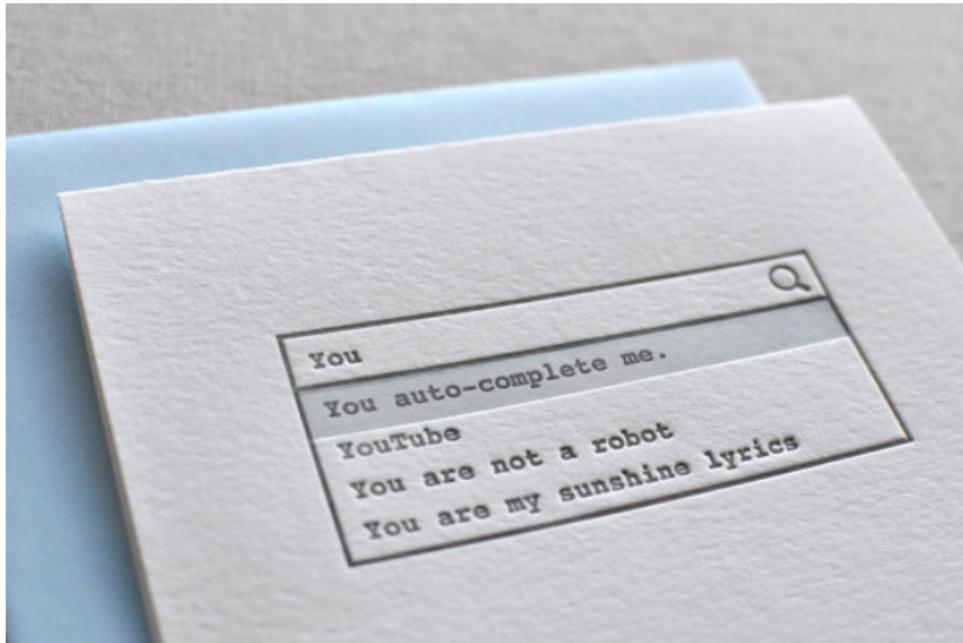
Evaluation



Evaluation



But we don't use topic models for prediction!



Qualitative Evaluation of the Latent Space

“segment 1”	“segment 2”	“matrix 1”	“matrix 2”	“line 1”	“line 2”	“power 1”	power 2”
imag SEGMENT texture color tissue brain slice cluster mri volume	speaker speech recogni signal train hmm source speakerind. SEGMENT sound	robust MATRIX eigenvalu uncertainti plane linear condition perturb root suffici	manufactur cell part MATRIX cellular famili design machinepart format group	constraint LINE match locat imag geometr impos segment fundament recogn	alpha redshift LINE galaxi quasar absorp high ssup densiti veloc	POWER spectrum omega mpc hsup larg redshift galaxi standard model	load memori vlsi POWER systolic input complex arra present implement

Figure 3: Eight selected factors from a 128 factor decomposition. The displayed word stems are the 10 most probable words in the class-conditional distribution $P(w|z)$, from top to bottom in descending order.

?

Qualitative Evaluation of the Latent Space

“Arts”	“Budgets”	“Children”	“Education”
NEW	MILLION	CHILDREN	SCHOOL
FILM	TAX	WOMEN	STUDENTS
SHOW	PROGRAM	PEOPLE	SCHOOLS
MUSIC	BUDGET	CHILD	EDUCATION
MOVIE	BILLION	YEARS	TEACHERS
PLAY	FEDERAL	FAMILIES	HIGH
MUSICAL	YEAR	WORK	PUBLIC
BEST	SPENDING	PARENTS	TEACHER
ACTOR	NEW	SAYS	BENNETT
FIRST	STATE	FAMILY	MANIGAT
YORK	PLAN	WELFARE	NAMPHY
OPERA	MONEY	MEN	STATE
THEATER	PROGRAMS	PERCENT	PRESIDENT
ACTRESS	GOVERNMENT	CARE	ELEMENTARY
LOVE	CONGRESS	LIFE	HAITI

?

Qualitative Evaluation of the Latent Space

- DA centralbank europæiske ecb s lån centralbanks
DE zentralbank ezb bank europäischen investitionsbank darlehen
EL τράπεζα τράπεζας κεντρική εκτ κεντρικής τράπεζες
EN **bank central ecb banks european monetary**
ES banco central europeo bce bancos centrales
FI keskuspankin ekp n europan keskuspankki eip
FR banque centrale bce européenne banques monétaire
IT banca centrale bce europea banche prestiti
NL bank centrale ecb europese banken leningen
PT banco central europeu bce bancos empréstimos
SV centralbanken europeiska ecb centralbankens s lån

?

Qualitative Evaluation of the Latent Space

(a) Topic labeled as SSL

Keyword	Probability
ssl	0.373722
expr	0.042501
init	0.033207
engine	0.026447
var	0.022222
ctx	0.023067
ptemp	0.017153
mctx	0.013773
lookup	0.012083
modssl	0.011238
ca	0.009548

(b) Topic labeled as Logging

Keyword	Probability
log	0.141733
request	.036017
mod	0.0311
config	0.029871
name	0.023725
headers	0.021266
autoindex	0.020037
format	0.017578
cmd	0.01512
header	0.013891
add	0.012661

Table 2: Sample Topics extracted from Apache source code

?

Word Intrusion

- Take the highest probability words from a topic

Original Topic

dog

cat

horse

pig

cow

Word Intrusion

- Take the highest probability words from a topic

Original Topic

dog

cat

apple

horse

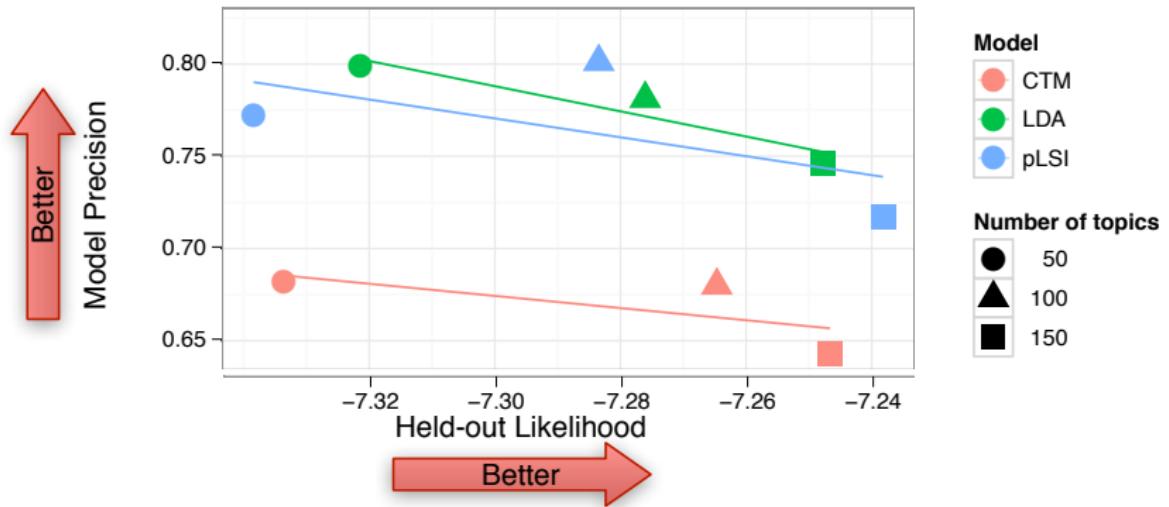
pig

cow

- Intruder: high probability word from another topic

Interpretability and Likelihood

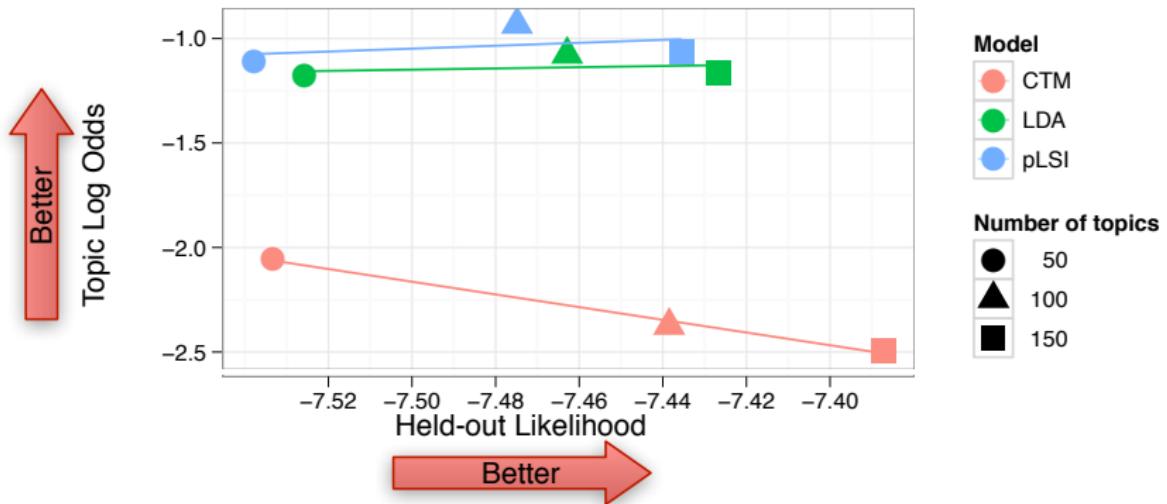
Model Precision on New York Times



within a model, higher likelihood \neq higher interpretability

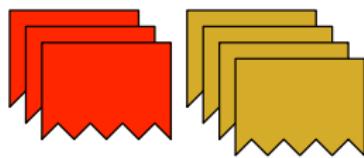
Interpretability and Likelihood

Topic Log Odds on Wikipedia



across models, higher likelihood \neq higher interpretability

What about Supervised Models?

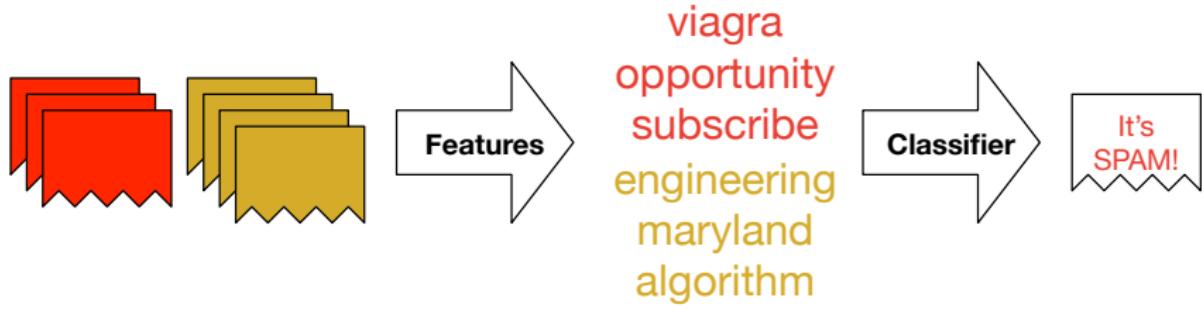


What about Supervised Models?

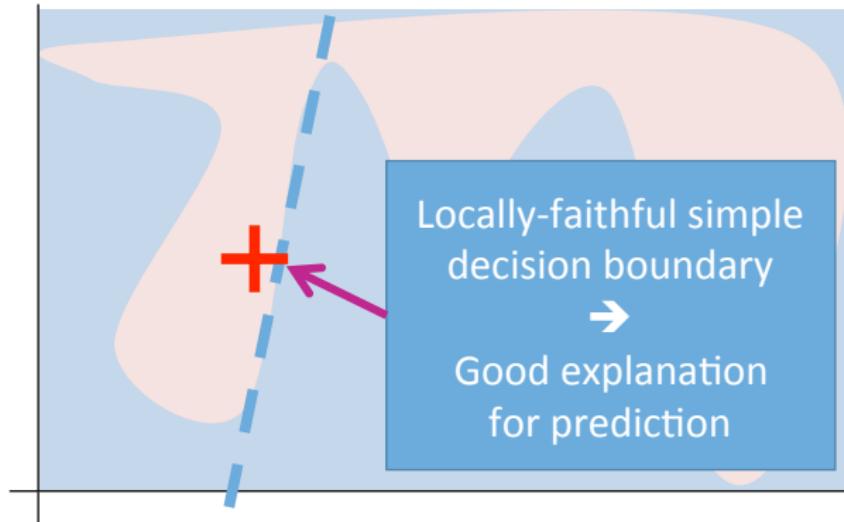


viagra
opportunity
subscribe
engineering
maryland
algorithm

What about Supervised Models?



LIME



Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "Why Should I Trust You?" Explaining the Predictions of Any Classifier. KDD 2016.
LIME: Local Interpretable Model-Agnostic Explanations

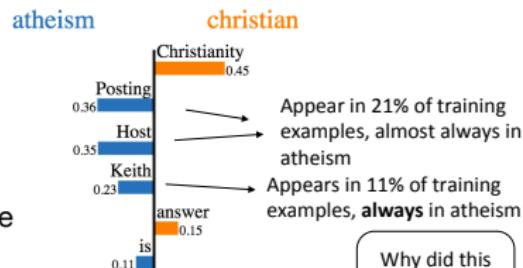
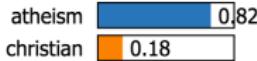
What's an Explanation

From: Keith Richards
Subject: Christianity is the answer
NTTP-Posting-Host: x.x.com

I think Christianity is the one true religion.
If you'd like to know more, send me a note



Prediction probabilities



Why did this happen? How do I fix it?



What's an Explanation



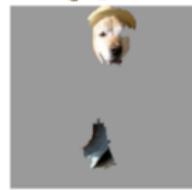
$P(\text{Guitarist}) = 0.32$



$P(\text{Guitarist}) = 0.24$



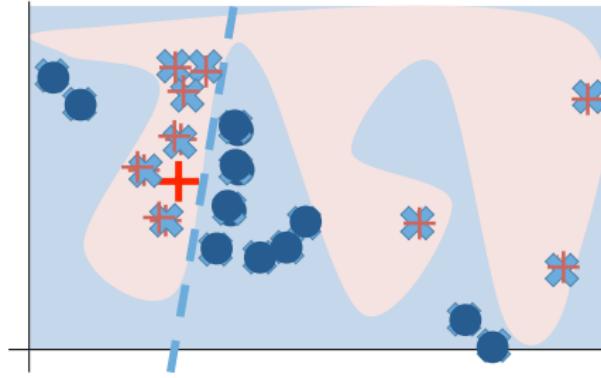
$P(\text{Guitarist}) = 0.21$



What makes good Explanation?

- Interpretable: Humans can Understand
- Faithful: Describes Model
- Model Agnostic: Generalize to Many Models

Method



- Complicated model predicts “near” example
- Simple model explains **local variation**
- **Explains what complicated model focused on**

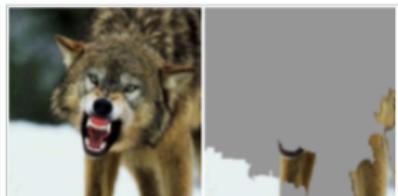
Is this a good Classifier?



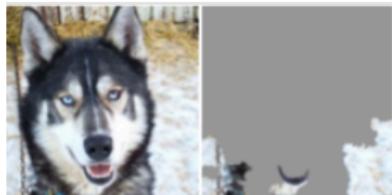
Predicted: **wolf**
True: **wolf**



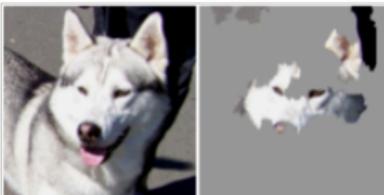
Predicted: **husky**
True: **husky**



Predicted: **wolf**
True: **wolf**



Predicted: **wolf**
True: **husky**

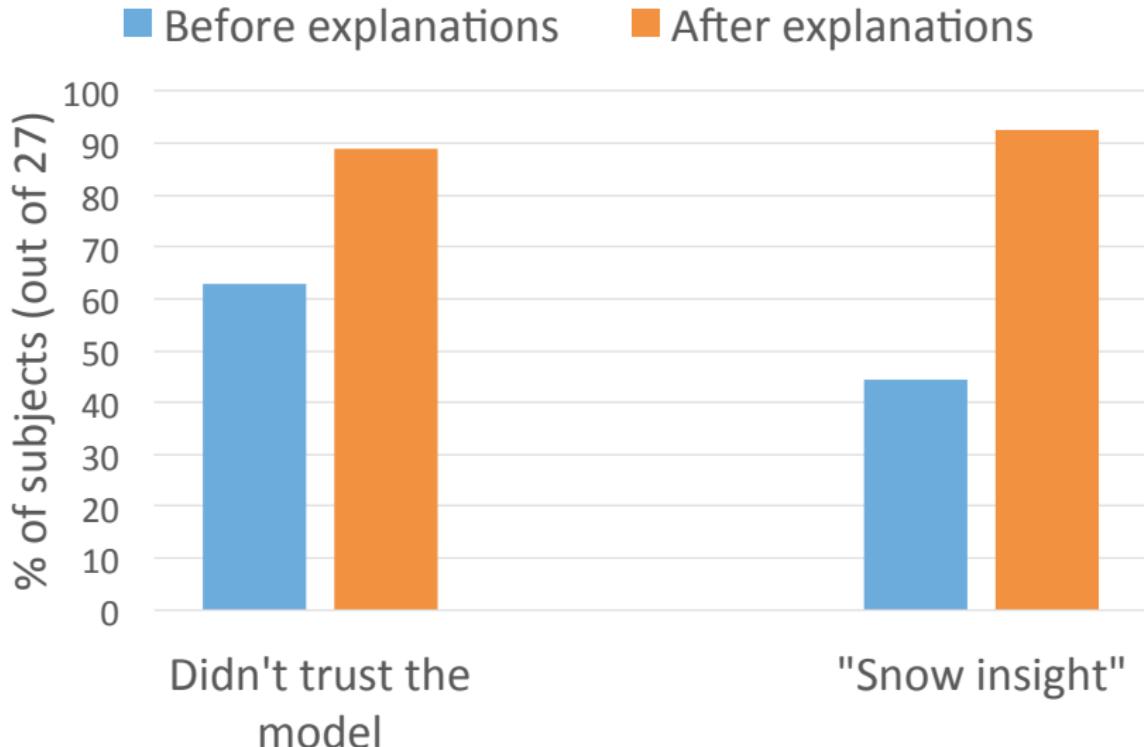


Predicted: **husky**
True: **husky**

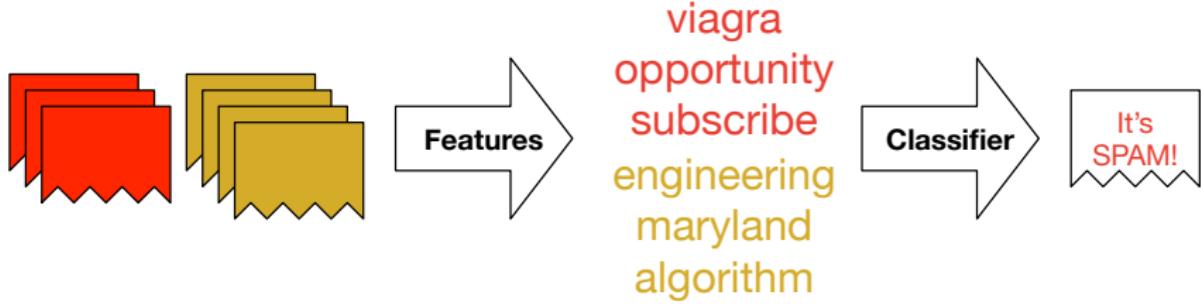


Predicted: **wolf**
True: **wolf**

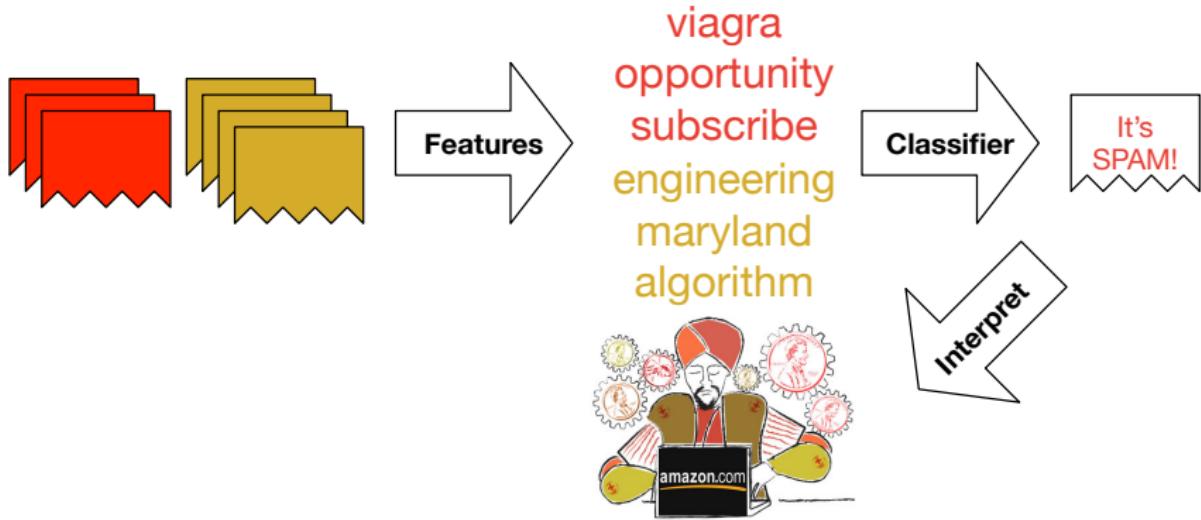
Is this a good Classifier?



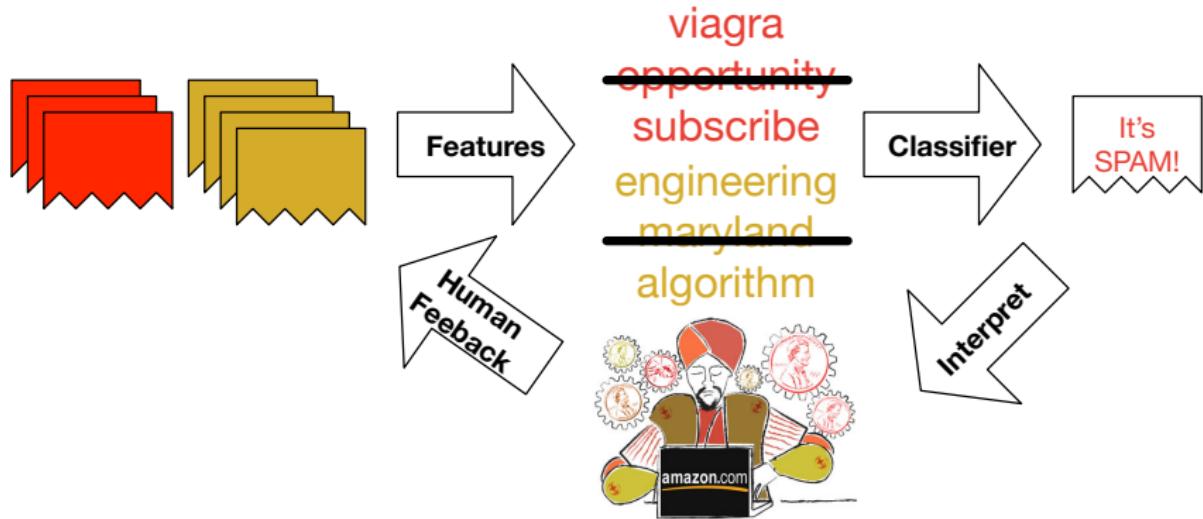
Improving ML Algorithms



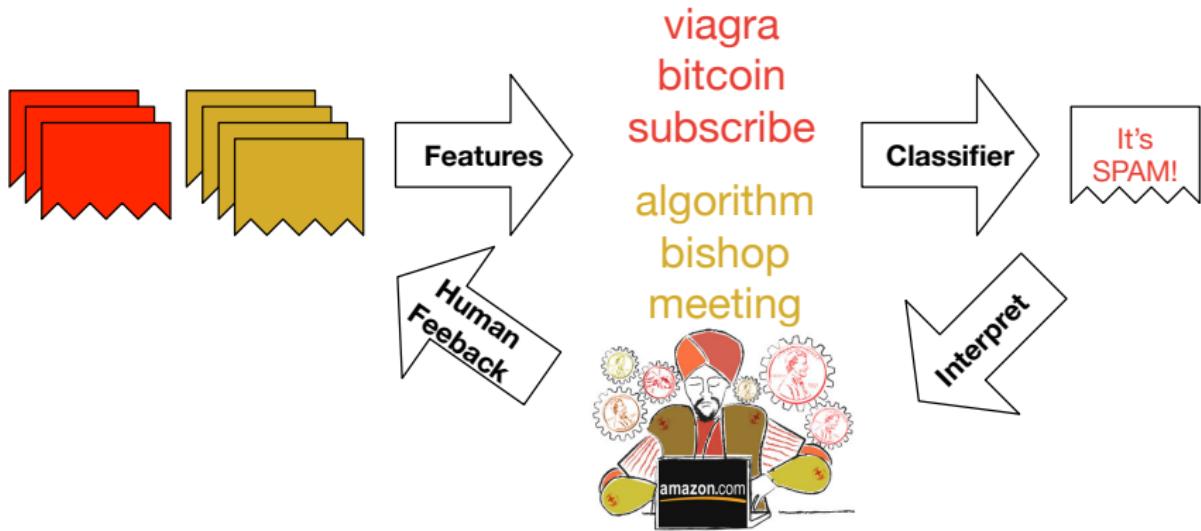
Improving ML Algorithms



Improving ML Algorithms



Improving ML Algorithms



Improving ML Algorithms

Example #5 of 10

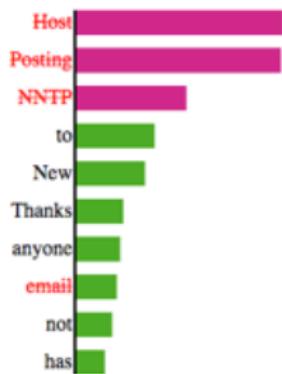
True Class:  Atheism

[Instructions](#)

[Previous](#)

[Next](#)

Words that the algorithm considers important.



Bar length indicates importance, and color indicates to which topic: Christianity (green) or Atheism (Pink).

Please click on the words (right next to the bars) that you think the algorithm is using incorrectly, because they are not important to distinguish between Atheism and Christianity. They should be red and crossed off after you click them.

Document

From: johnchad@triton.unm.edu (jchadwic)

Subject: Another request for Darwin Fish

Organization: University of New Mexico, Albuquerque

Lines: 11

NNTP-Posting-Host: triton.unm.edu

Hello Gang,

There have been some notes recently asking where **to** obtain the DARWIN fish.

This is the same question I have and I have **not** seen an answer on the net. If **anyone has** a contact please post on the net or **email** me.

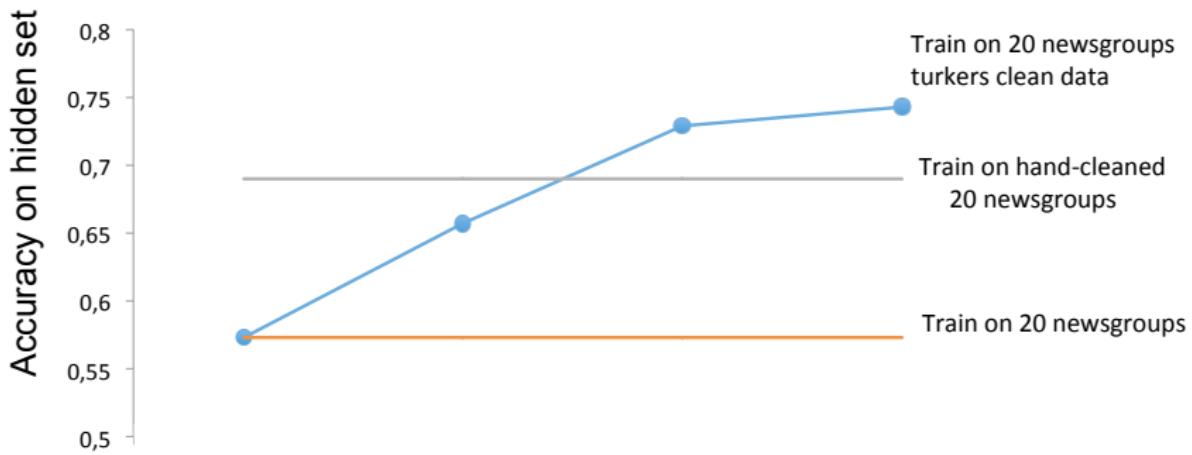
Thanks,

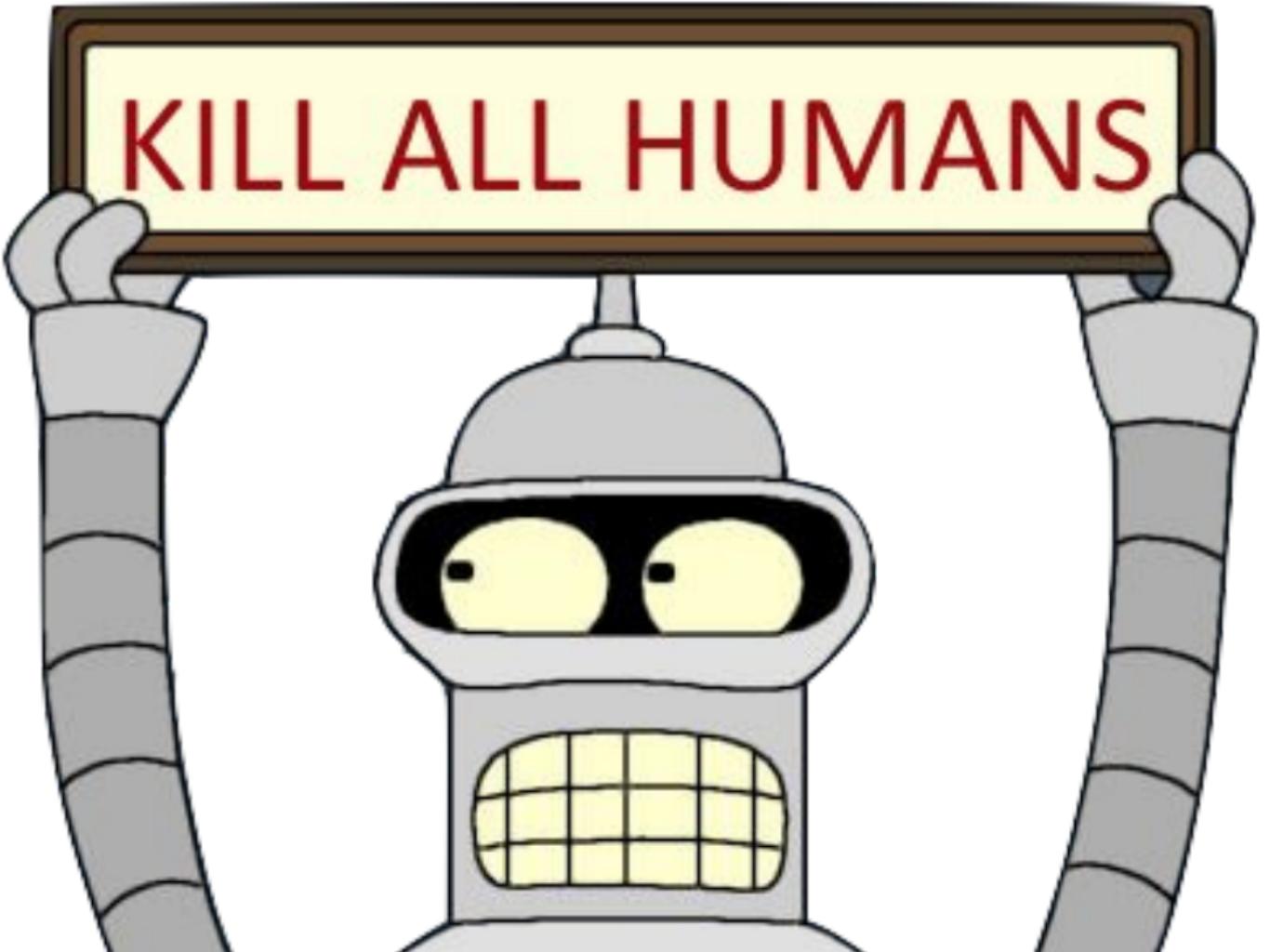
john chadwick

johnchad@triton.unm.edu

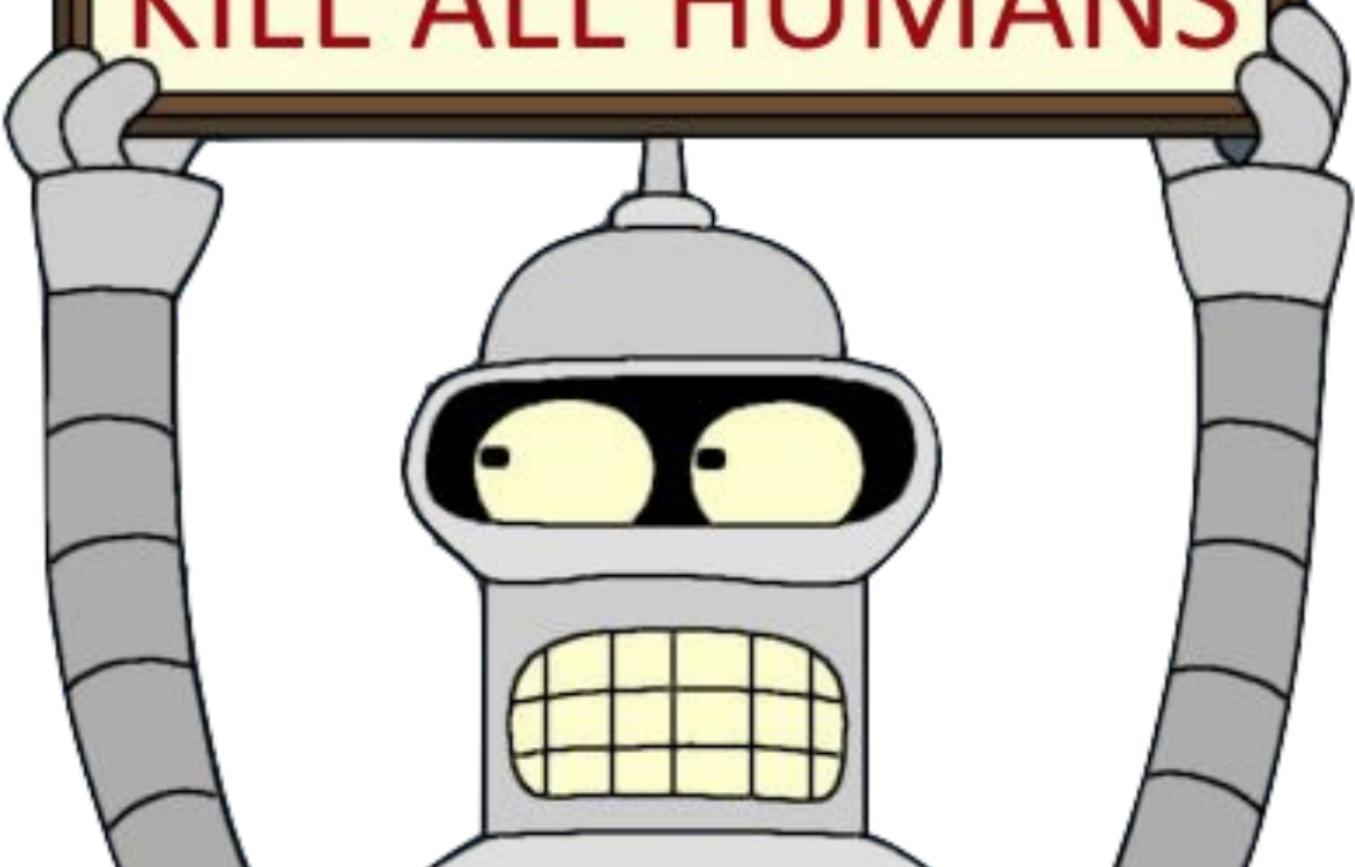
or

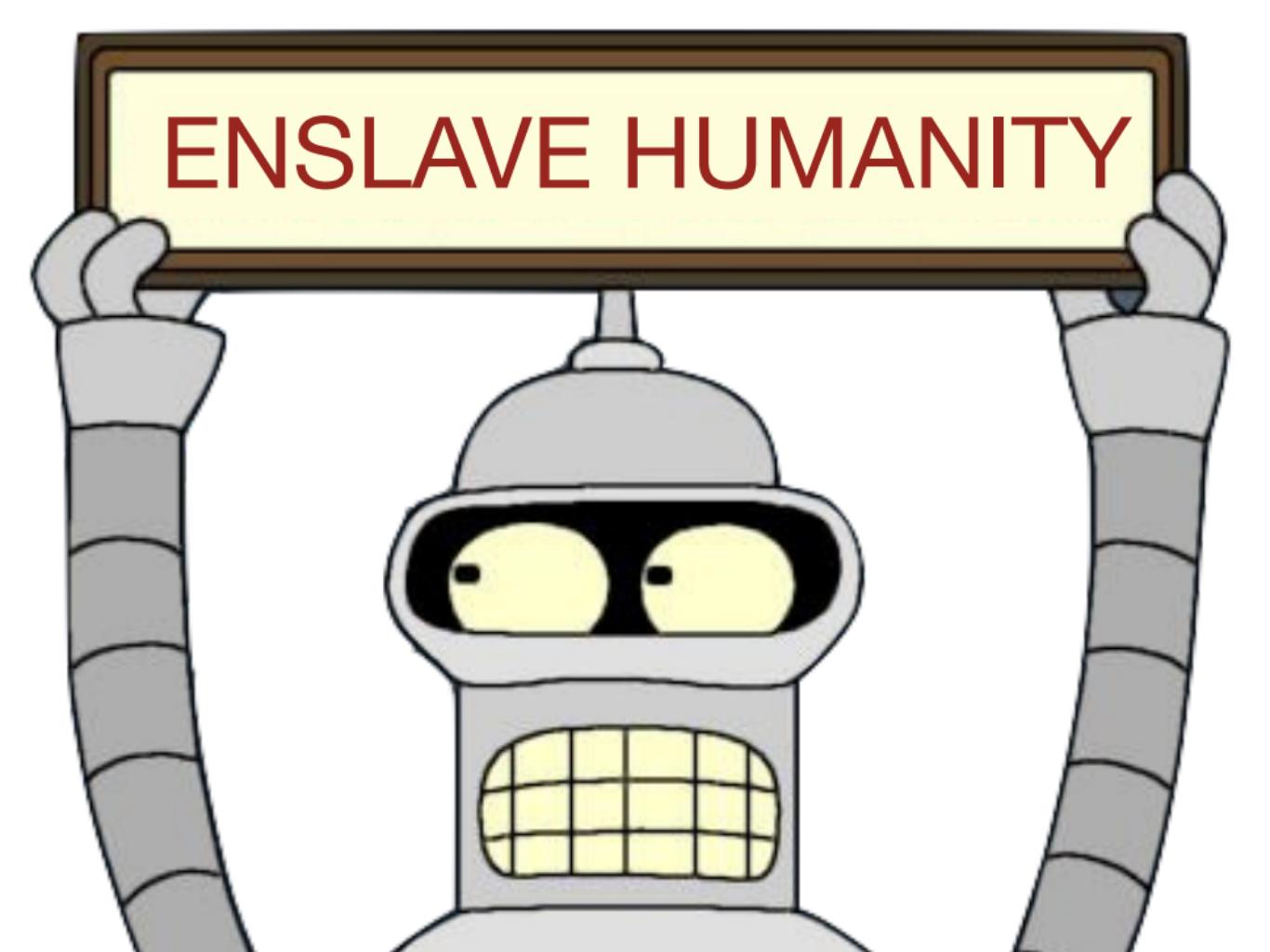
Improving ML Algorithms





KILL ALL HUMANS



A cartoon illustration of a silver robot head with large, white, almond-shaped eyes and a mouth made of a grid of yellow squares. The robot is holding a rectangular sign with a brown border and a light yellow background. The sign features the text "ENSLAVE HUMANITY" in bold, red, sans-serif capital letters.

ENSLAVE HUMANITY





початок об 11.00

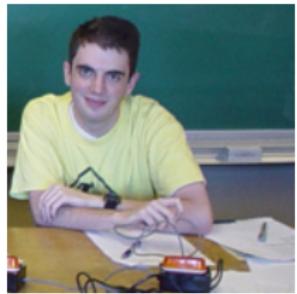
Готель Radisson Blu
м. Київ, проспект Голославів Вал 22

Олена
БОЙЧУН



Centaur Chess

Measuring Interpretability

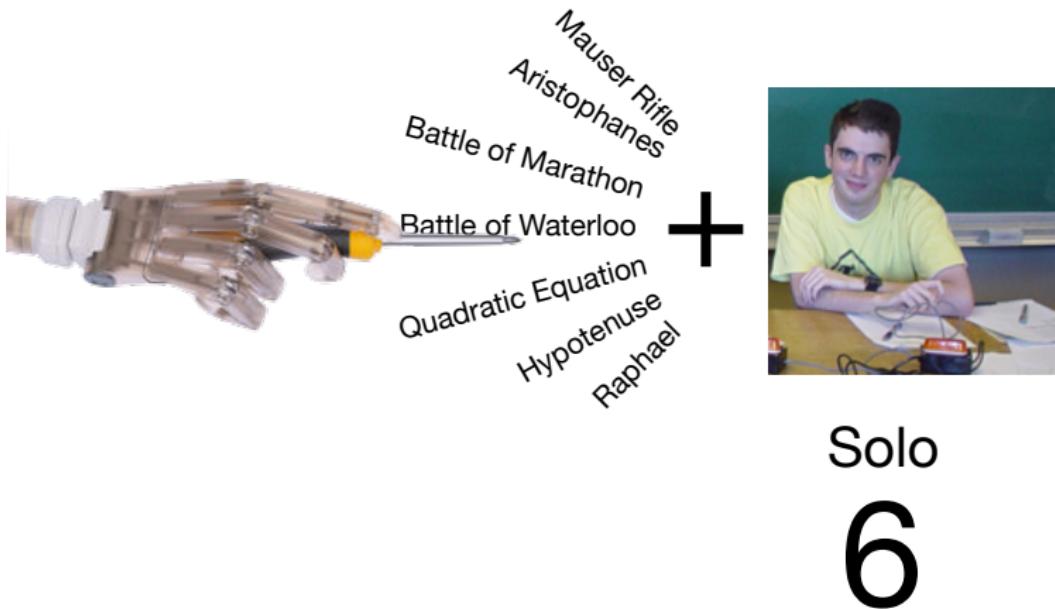


Measuring Interpretability

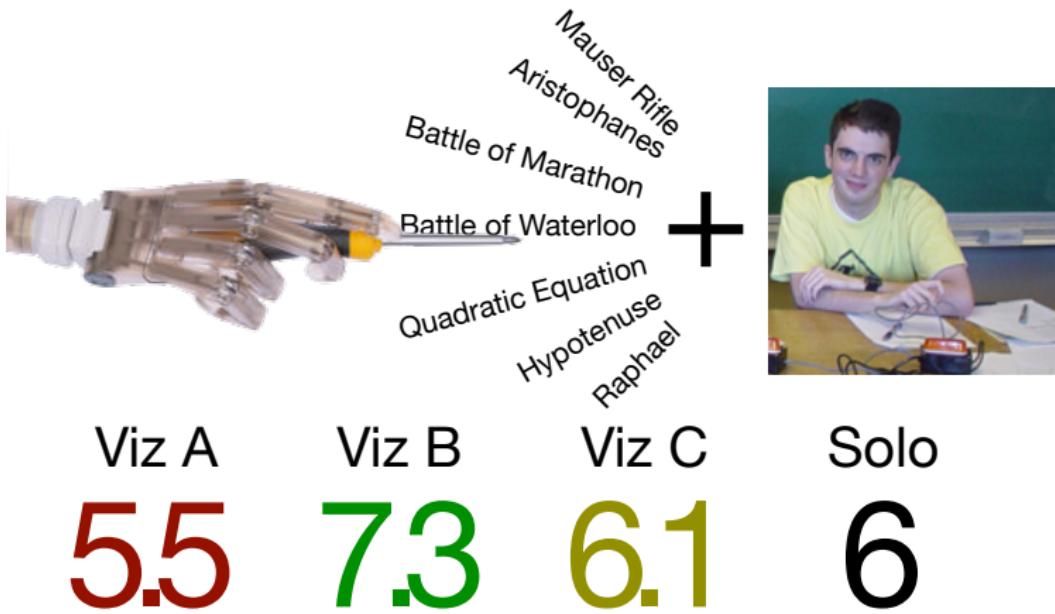


Solo
6

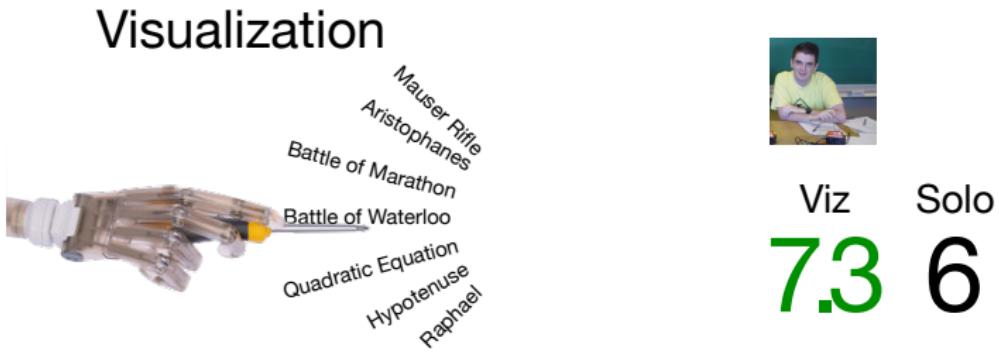
Measuring Interpretability



Measuring Interpretability



Improvement through Reinforcement Learning



Improvement through Reinforcement Learning



Improvement through Reinforcement Learning



Improvement through Reinforcement Learning



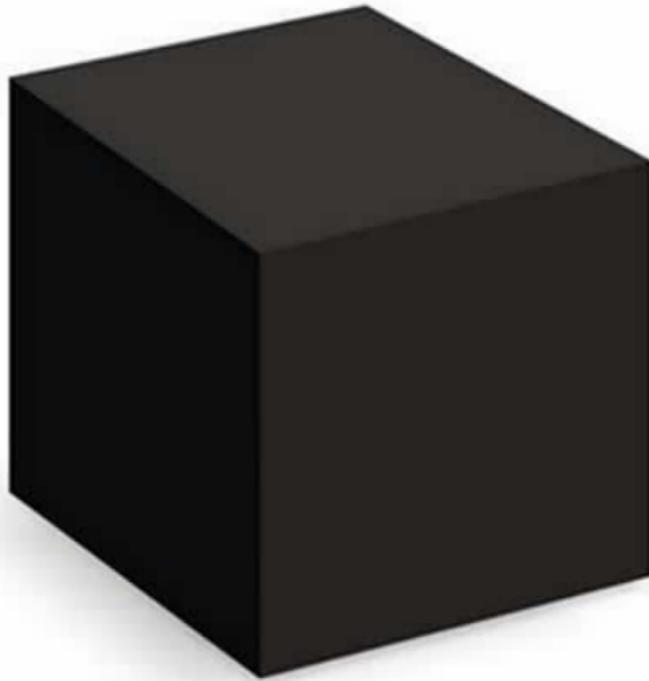
Improvement through Reinforcement Learning

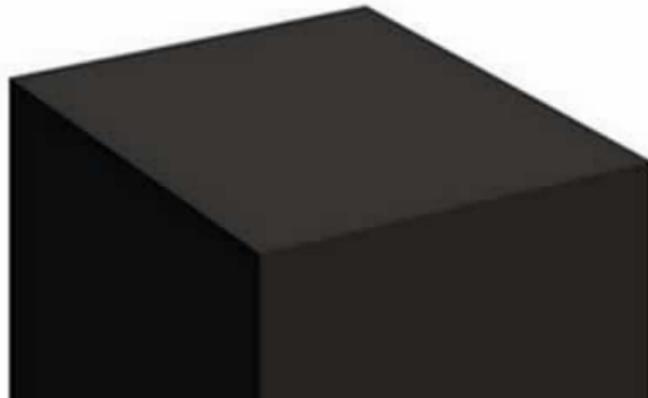


Simultaneous Interpretation is Hard!

- Exhausting for humans
- Computers not trusted
- Differential strengths
- Same word-by-word characteristic







Takeaways

- ML should be interpretable
- We should measure interpretability
- Interpretability should reflect the world we want