

Human Ranking of Machine Translation

Matt Post
Johns Hopkins University

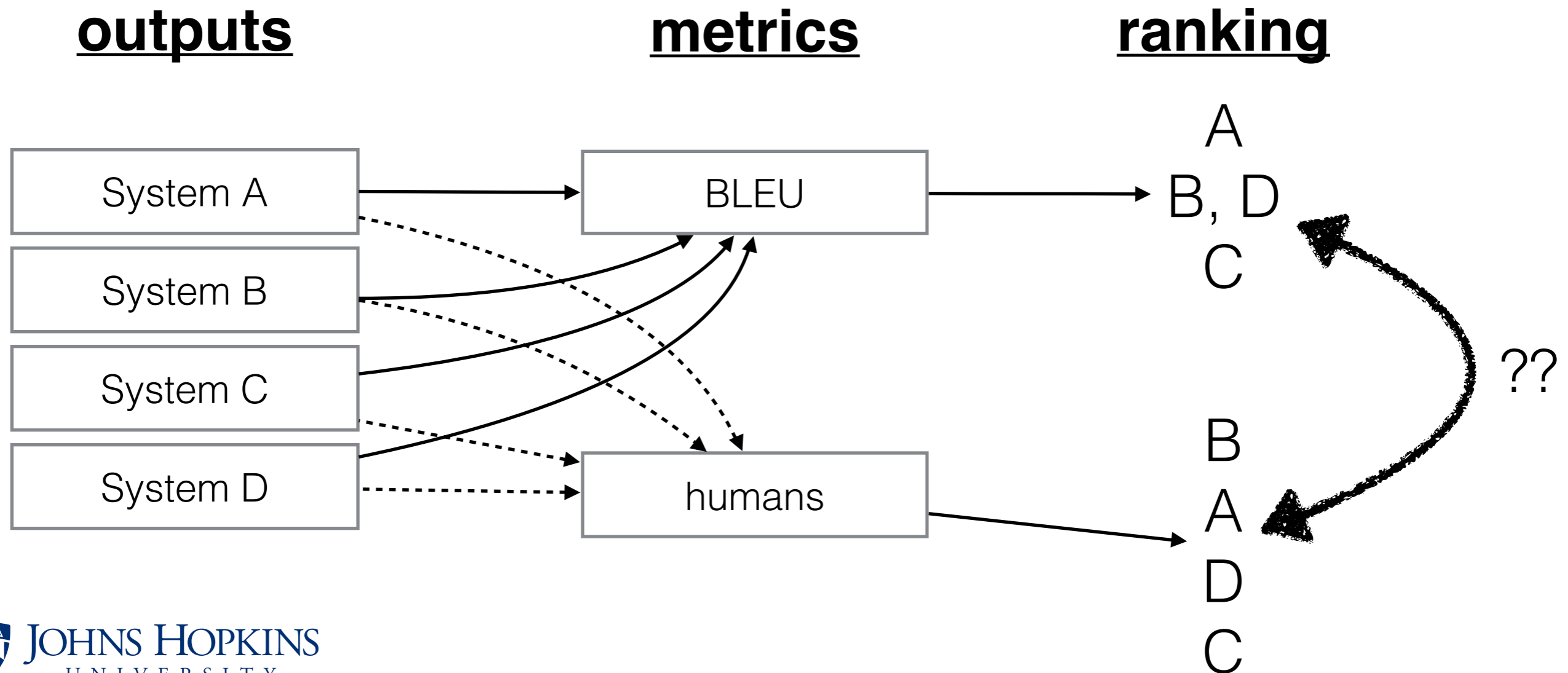
University of Pennsylvania
April 9, 2015

Review

- In translation, human evaluations are what matter
 - but they are expensive to run
 - this holds up science!
- The solution is automatic metrics
 - fast, cheap, (usually) easy to compute
 - deterministic

Review

- Automatic metrics produce a *ranking*
- They are evaluated using correlation statistics against *human judgments*



Review

- The human judgments are the “gold standard”
- Questions:
 1. How do we get this gold standard?
 2. How do we know it's correct?



Today

- **How we produce the gold-standard ranking**
- **How we know it's correct**

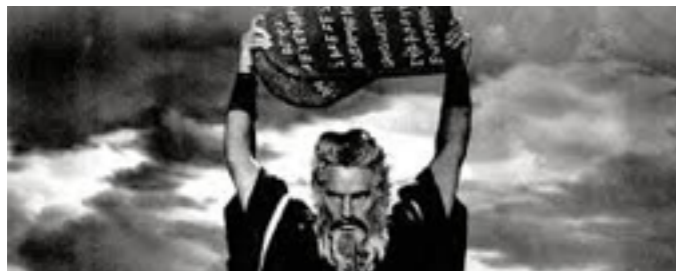
At the end of this lecture...

- You should understand
 - how to rank with incomplete
 - how to evaluate truth claims in science
- You might come away with
 - a desire to submit your metric to the WMT metrics task (deadline: May 25, 2015)
 - a desire to buy an Xbox
 - a preference for simplicity

Producing a ranking

- Then, we take this data and produce a ranking
- Outline of the rest of the talk

Human ranking methods



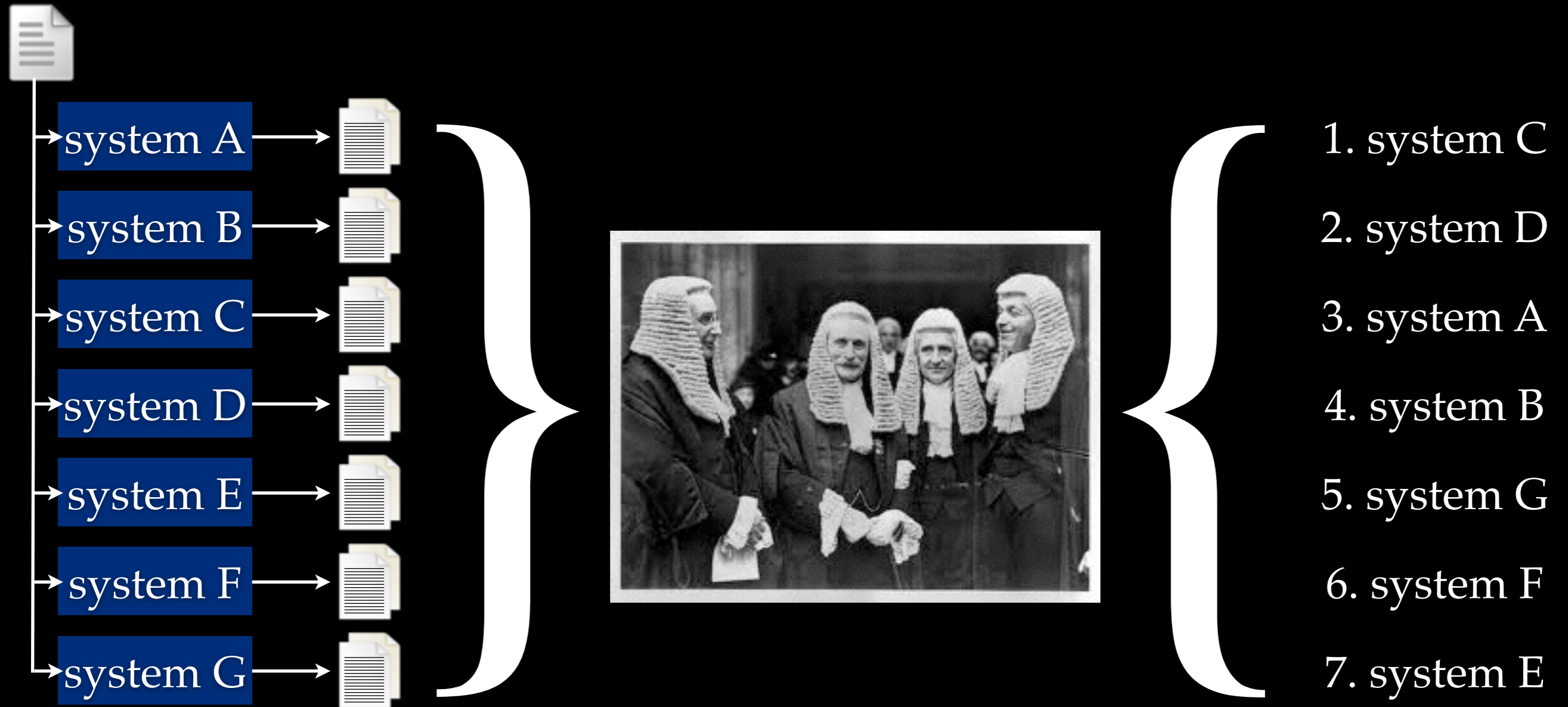
Model selection



Clustering



Goal



Goal

- Produce a ranking of *systems*
- There are many ways to do this:
 - Reading comprehension tests
 - Time spent on human post-editing
 - Aggregating sentence-level judgments
- This last one is what is used by the Workshop on Statistical Machine Translation (statmt.org/wmt15)

Inherent problems

- Translation is used for a range of tasks



Understanding the past



Technical manuals



Conversing



Information

- What *best* (or sufficient) means likely varies by person and situation

Collecting data

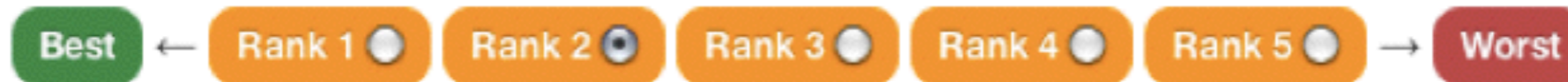
- Data: K systems translate an N-sentence document
- We use human judges to compare translations of an input sentence and select whether
the first is *better*,
worse, or
equivalent to the second
- We use a large pool of judges

"Valentino měl vždycky raději eleganci než slávu.

— Source

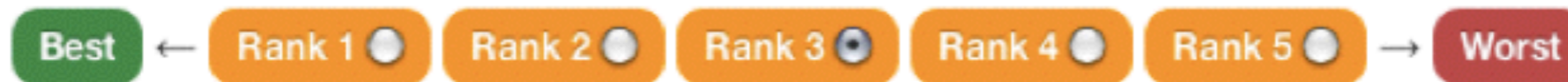
Valentino has always preferred elegance to notoriety.

— Reference



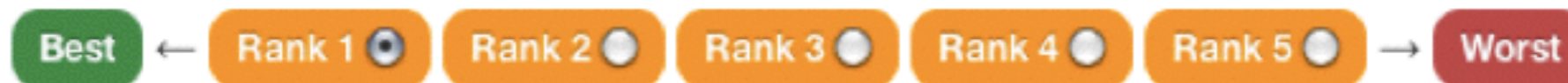
"Valentino should always elegance rather than fame.

— Translation 1



"Valentino has always rather than the elegance of glory.

— Translation 2



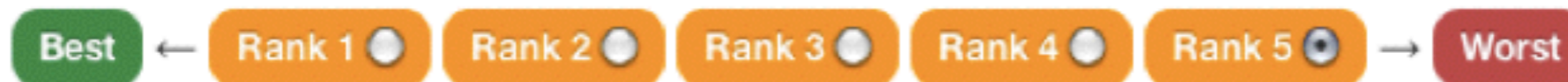
" Valentino had always preferred elegance than glory.

— Translation 3



"Valentino has always had the elegance rather than glory.

— Translation 4



" Valentino has always had a rather than the elegance of the glory.

— Translation 5

Collecting data

"Valentino měl vždycky raději eleganci než slávu.

— Source

Valentino has always preferred elegance to notoriety.

— Reference

Best ← Rank 1 Rank 2 Rank 3 Rank 4 Rank 5 → Worst

"Valentino should always elegance rather than fame.

— Translation 1

Best ← Rank 1 Rank 2 Rank 3 Rank 4 Rank 5 → Worst

"Valentino has always rather than the elegance of glory.

— Translation 2

Best ← Rank 1 Rank 2 Rank 3 Rank 4 Rank 5 → Worst

"Valentino had always preferred elegance than glory.

— Translation 3

Best ← Rank 1 Rank 2 Rank 3 Rank 4 Rank 5 → Worst

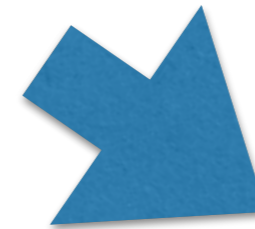
"Valentino has always had the elegance rather than glory.

— Translation 4

Best ← Rank 1 Rank 2 Rank 3 Rank 4 Rank 5 → Worst

"Valentino has always had a rather than the elegance of the glory.

— Translation 5

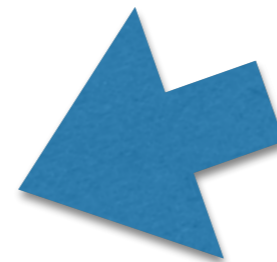


$C > A > B > D > E$



$C > A$	$A > B$	$B > D$	$D > E$
$C > B$	$A > D$	$B > E$	
$C > D$	$A > E$		
$C > E$			

ten pairwise judgments



Dataset

- This yields ternary-valued pairwise judgments of the following form

judge “dredd” ranked onlineB $>$ JHU on sent #74

judge “judy” ranked uedin $>$ UU on sent #1734

judge “reinhold” ranked JHU $>$ UU on sent #1

judge “jay” ranked onlineA = uedin on sent #953

...

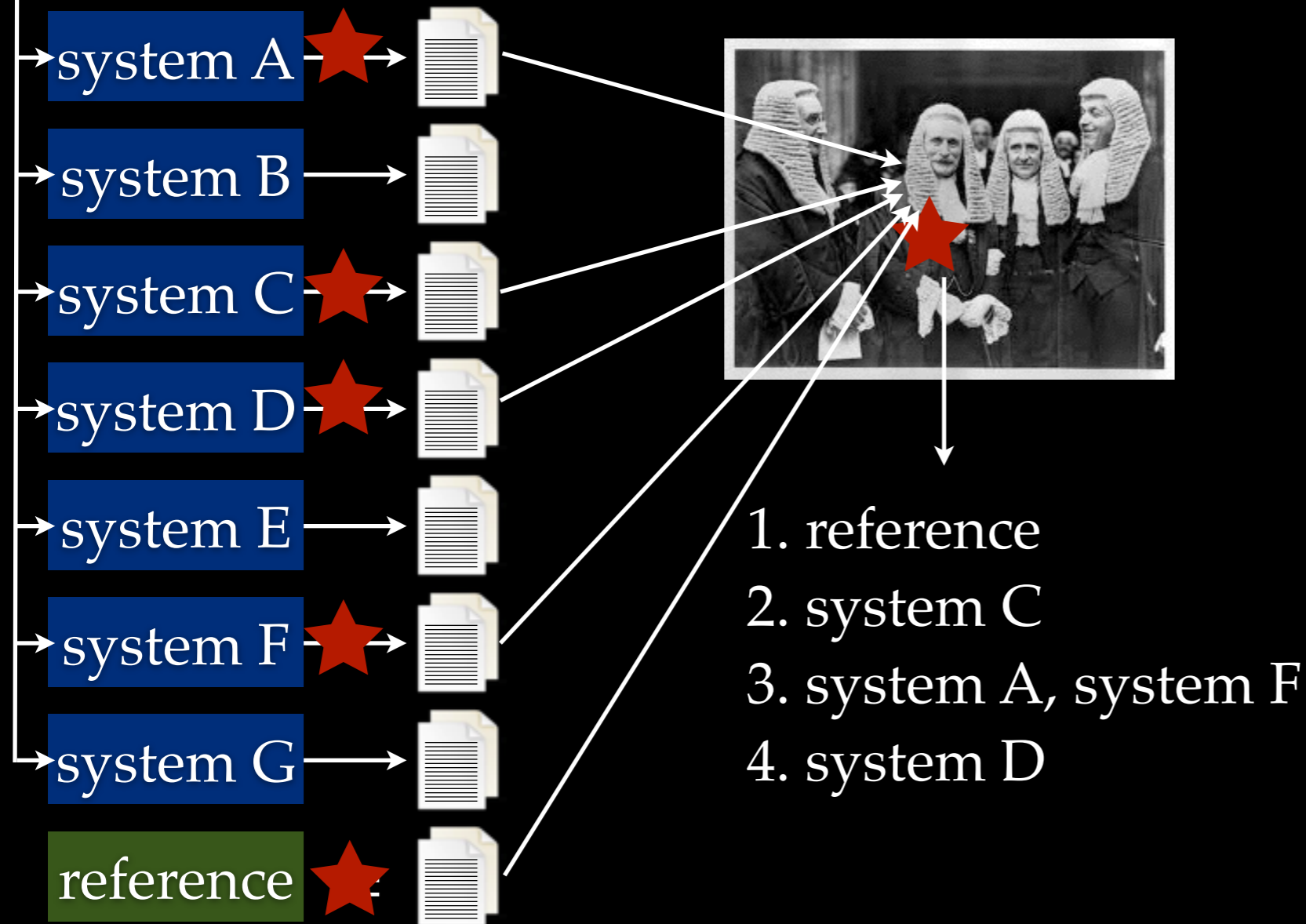
The sample space

- How much data is there to collect?

(number of ways to pick two systems)
x (number of sentences) x (number of judges)

- For 10 systems there are 135k comparisons
 - For 20 systems, 570k
 - More with multiple judges
- Too much to collect, also wasteful; instead we sample

Design of the WMT Evaluation (2008-2011)



WMT Raw Data:
pairwise rankings

- reference \prec system A
- reference \prec system C
- reference \prec system D
- reference \prec system F
- system A \succ system C
- system A \prec system D
- system A \equiv system F
- system C \prec system D
- system C \prec system F
- system D \prec system F

While (evaluation period is not over):

- Sample input sentence.
- Sample five translators of it from $Systems \cup \{Reference\}$.
- Sample a judge.
- Receive set of pairwise judgments from the judge.

How much data do we collect?

LANGUAGE PAIR	Systems	Rankings	Average
Czech–English	5	21,130	4,226.0
English–Czech	10	55,900	5,590.0
German–English	13	25,260	1,943.0
English–German	18	54,660	3,036.6
French–English	8	26,090	3,261.2
English–French	13	33,350	2,565.3
Russian–English	13	34,460	2,650.7
English–Russian	9	28,960	3,217.7
Hindi–English	9	20,900	2,322.2
English–Hindi	12	28,120	2,343.3
TOTAL WMT 14	110	328,830	2,989.3
WMT13	148	942,840	6,370.5
WMT12	103	101,969	999.6
WMT11	133	63,045	474.0

of tens of
millions
possible

Producing a ranking

- Then, we take this data and produce a ranking
- Human ranking methods



Expected wins and variants



Bayesian model (relative ability)



TrueSkill™

Expected wins (1)

- This most appealing and intuitive approach
- Define $wins(A)$, $ties(A)$, and $loses(A)$ as the number of times system A won, tied, or lost

- Score each system as follows

$$\text{score}(A) = \frac{wins(A) + ties(A)}{wins(A) + ties(A) + loses(A)}$$

- Now sort by scores

Expected wins (2)

- Do you see any problems with this?

$$\text{score}(A) = \frac{\text{wins}(A) + \text{ties}(A)}{\text{wins}(A) + \text{ties}(A) + \text{loses}(A)}$$

- Look at a judgments:

jud	one winner, one loser	onlineB > JHU on sent #74
jud	one winner, one loser	edin > UU on sent #1734
jud	one winner, one loser	ed JHU > UU on sent #1
jud	two winners, no losers	lineA = uedin on sent #953

Expected wins (3)

- A system is rewarded as much for a tie as for a win
 - ...and most systems are variations of the same underlying architecture, data



MOSES
statistical
machine translation
system

- New formula: throw away ties

$$\text{score}(A) = \frac{\text{wins}(A)}{\text{wins}(A) + \text{loses}(A)}$$

- Wait: Is this better?

A Grain of Salt for the WMT Manual Evaluation (Bojar et al., 2012)

Expected wins (4)

- Problem 2: the luck of the draw

$$\text{score}(A) = \frac{\text{wins}(A)}{\text{wins}(A) + \text{loses}(A)}$$

aggregation over
different sets of inputs
different competitors
different judges

- Consider a case where in reality $B > C$, but
 - B gets compared to a bunch of good systems
 - C gets compared to a bunch of bad systems
 - we could get $\text{score}(C) > \text{score}(B)$

Expected wins (5)

- This can happen!
 - Systems include a human reference translation
 - Also include really good unconstrained commercial systems

Expected wins (6)

- Even more problems:
 - remember that the scores for a system is the percentage of time it won in comparisons across all systems
 - what if $\text{score}(B) > \text{score}(C)$, but in direct comparisons, C was almost always better than B?
 - this leads to cycles in the ranking
- Is this a problem?



Summary

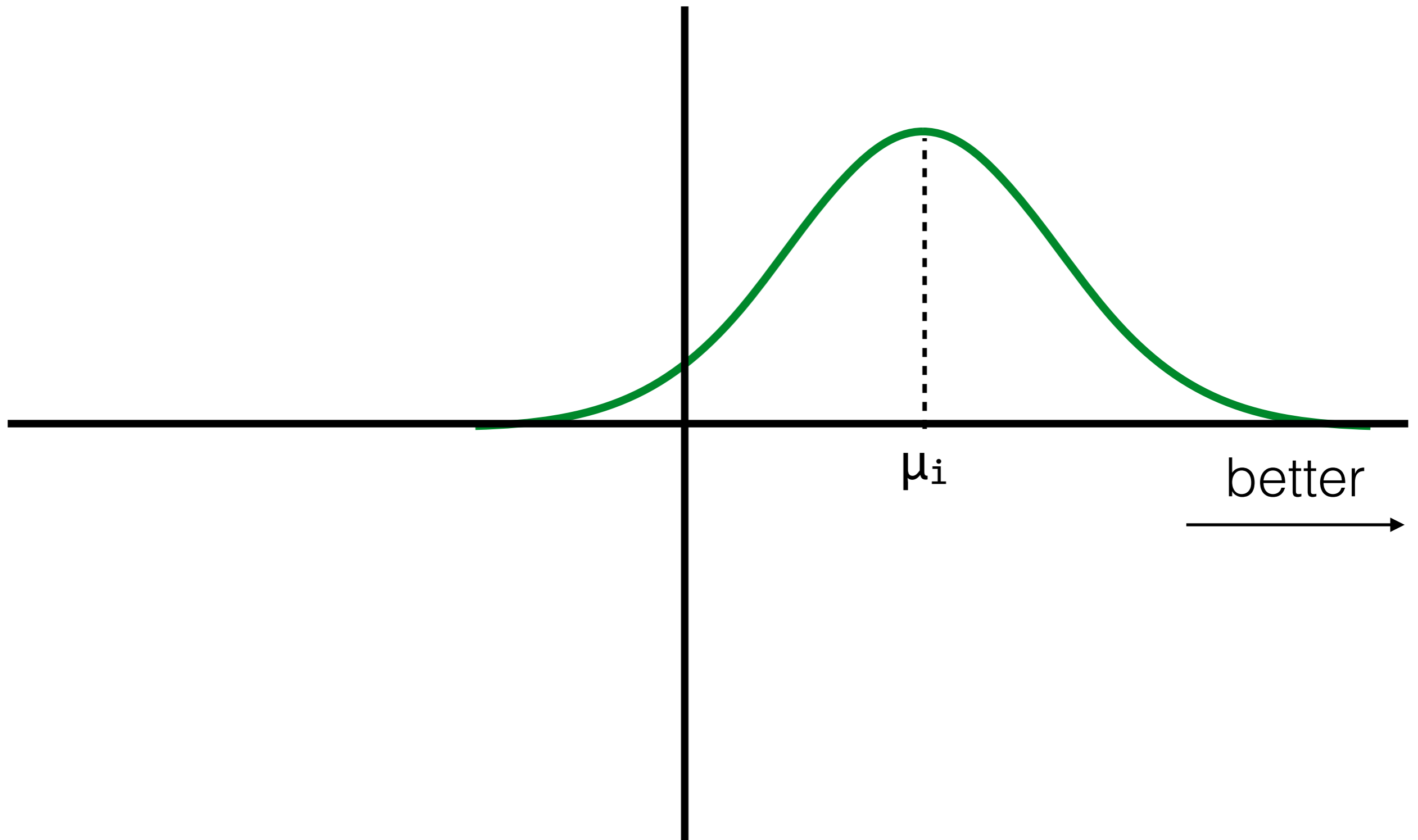
- List of problems:
 - Including ties biases similar systems, excluding discredits
 - Comparisons do not factor in difficulty of the “match” (i.e., losing to the best system should count less)
 - There are cycles in the judgments
- We made intuitive changes, but how do we know whether they’re correct?

Relative ability model

Models of Translation Competitions (Hopkins & May, 2013)

- In Expected Wins, we estimate a probability of each system winning a competition
- We now move to a setup that models the *relative ability* of a system
 - Assume each system S_i has an inherent ability, μ_j
 - Its translations are then represented by draws from a Gaussian distribution centered at μ_j

Relative ability



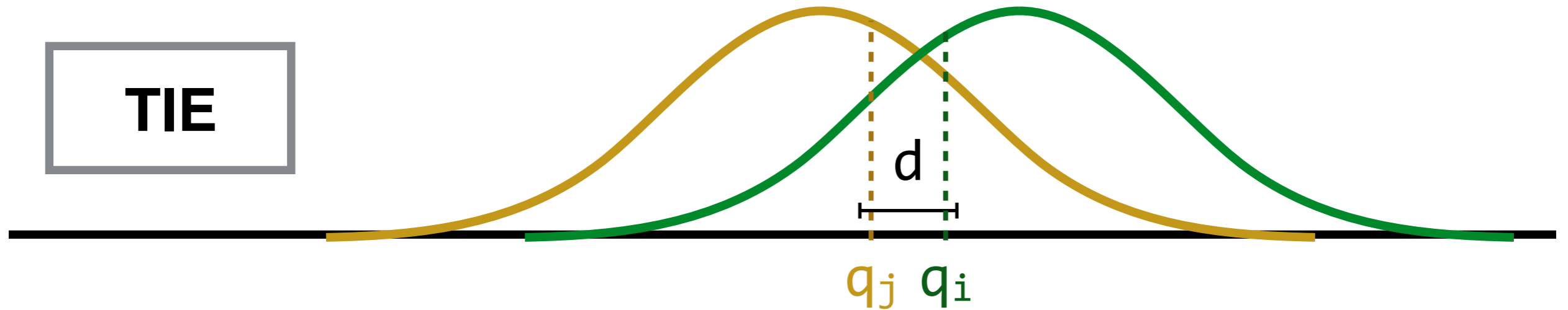
Relative ability

- A “competition” proceeds as follows:
 - Choose two systems, S_i and S_j , from the set $\{S\}$
 - Sample a “translation” from their distributions
$$q_i \sim N(S_i; \mu_i, \sigma^2)$$
$$q_j \sim N(S_j; \mu_j, \sigma^2)$$
 - Compare their values to determine who won
 - Define d as a “decision radius”
 - Record a tie if $|q_i - q_j| < d$
 - Else record a win or loss

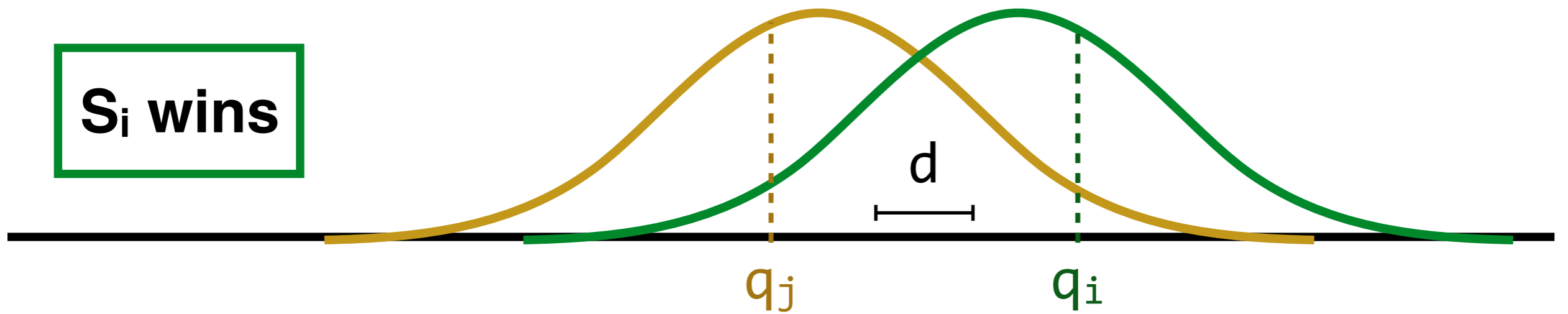
Visually

better \rightarrow

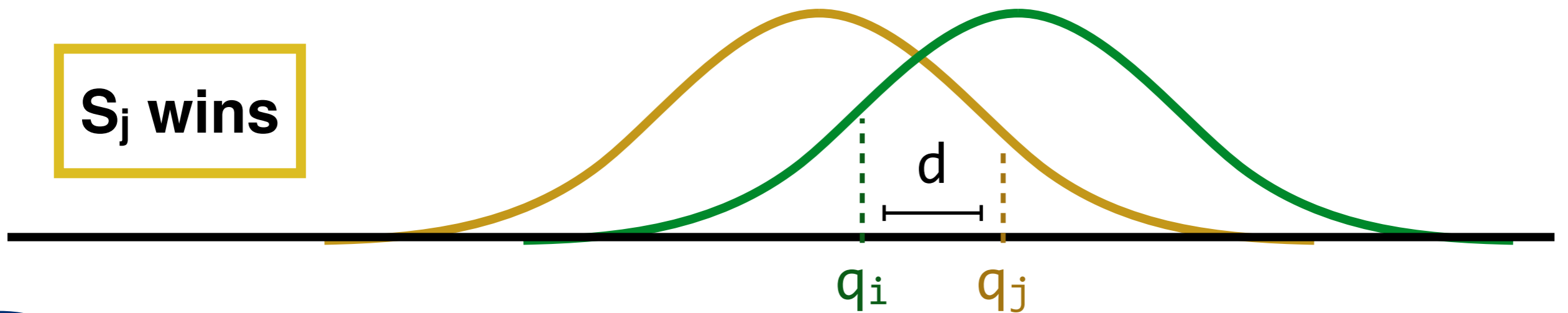
TIE



S_i wins



S_j wins



Observations

- We can compute exact probabilities for all these events (difference of Gaussians)
- On average, a system with a higher “ability” will have higher draws, and will win
- Systems with close μ s will tie more often

Learning the model

- If we knew the system means, we could rank them
- We assume the data was generated by the process above; we need to infer values for hidden params:
 - System means $\{\mu\}$
 - Sampled translation qualities $\{q\}$
- We'll use Gibbs sampling
 - Uses simple random steps to learn a complicated joint distributions
 - Converges under certain conditions

Gibbs sampling

judge “dredd” ranked onlineB > JHU on sent #74
judge “judy” ranked uedin > UU on sent #1734
judge “reinhold” ranked JHU > UU on sent #1
judge “jay” ranked onlineA = uedin on sent #953

- Represent data as tuples $(S_i, S_j, \pi, q_i, q_j)$

(onlineB, JHU, >, ?, ?)

(uedin, UU, >, ?, ?)

(JHU, UU, >, ?, ?)

(onlineA, uedin, =, ?, ?)

known

unknown

- Iterate back and forth between guessing $\{q\}$ s and $\{\mu\}$ s

Iterative process

```
[collect all the judgments]
until convergence
  # resample translation qualities
  for each judgment
     $q_i \sim N(\mu_i, \sigma^2)$ 
     $q_j \sim N(\mu_j, \sigma^2)$ 
    # (adjust samples to respect judgment  $\pi$ )

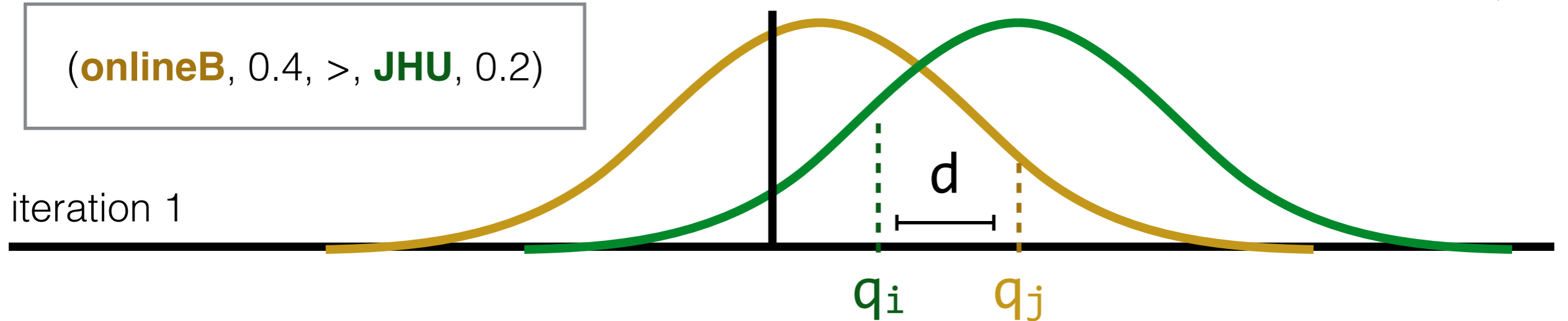
  # resample the system means
  for each system
     $\mu_i = \text{mean}(\{q_i\})$ 
```

Visually

better \rightarrow

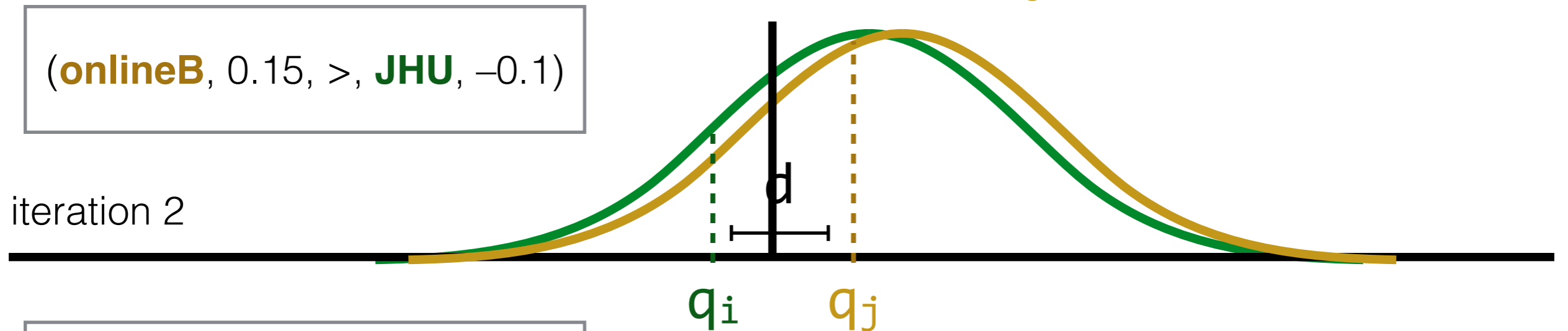
(**onlineB**, 0.4, >, **JHU**, 0.2)

iteration 1



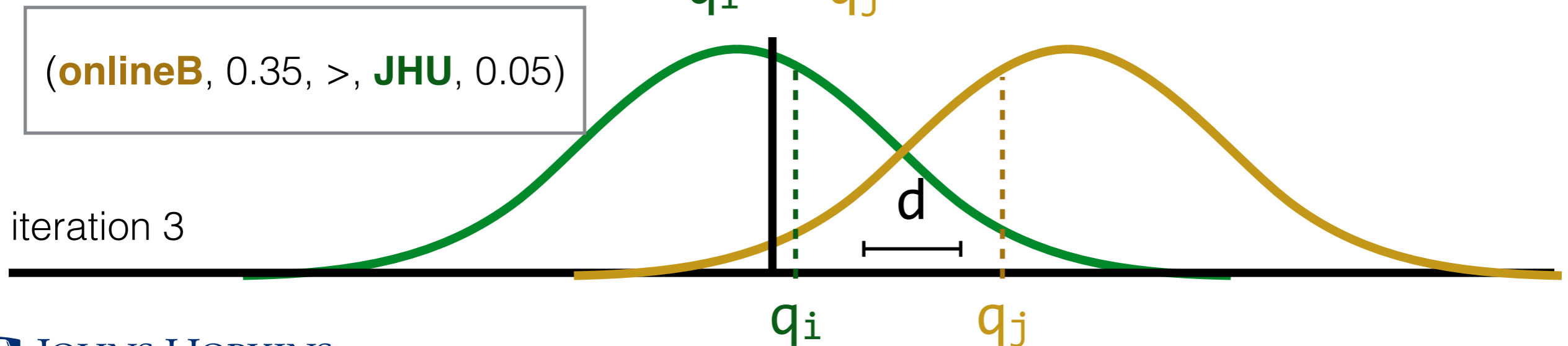
(**onlineB**, 0.15, >, **JHU**, -0.1)

iteration 2



(**onlineB**, 0.35, >, **JHU**, 0.05)

iteration 3



Summary

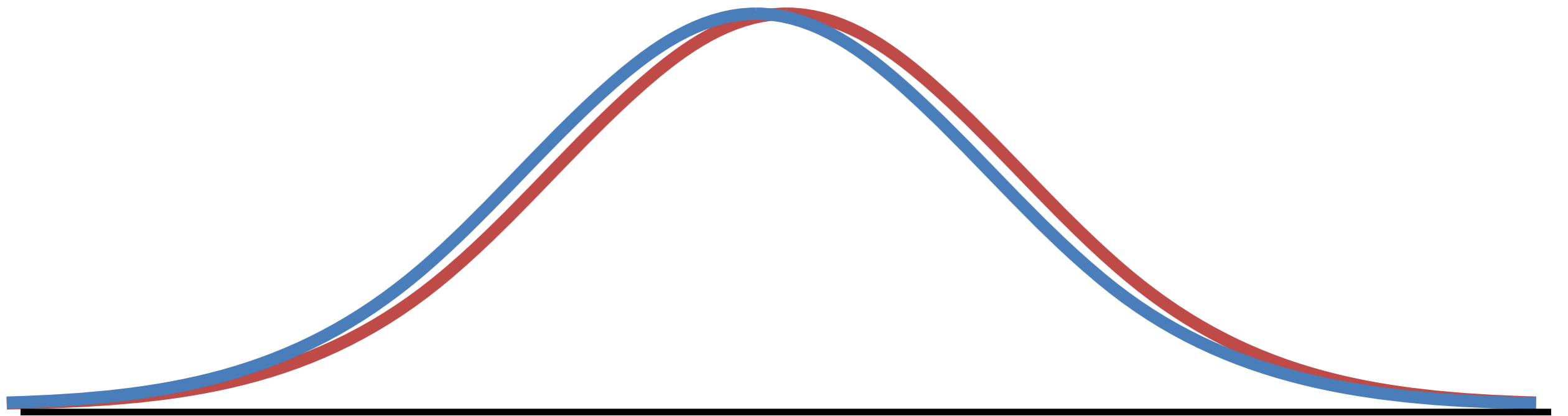
- Summary
 - Model provides us with an explanation of how the data was generated
 - We infer the abilities of the systems to rank using the human judgments
- Problems
 - Still no notion of evenness of the match
 - Judges are not modeled
 - Actual sentences are ignored

TrueSkill™ Ranking System

- Used to rate players in Xbox Live
- Based on the ELO system for Chess
- Models player ability (μ) *and* the system's confidence about that estimate (σ)
 - When a game is played, the outcome (win, loss, or tie) is used to update these parameters
 - A more surprising outcome results in larger updates
 - These values are also used to find even matches

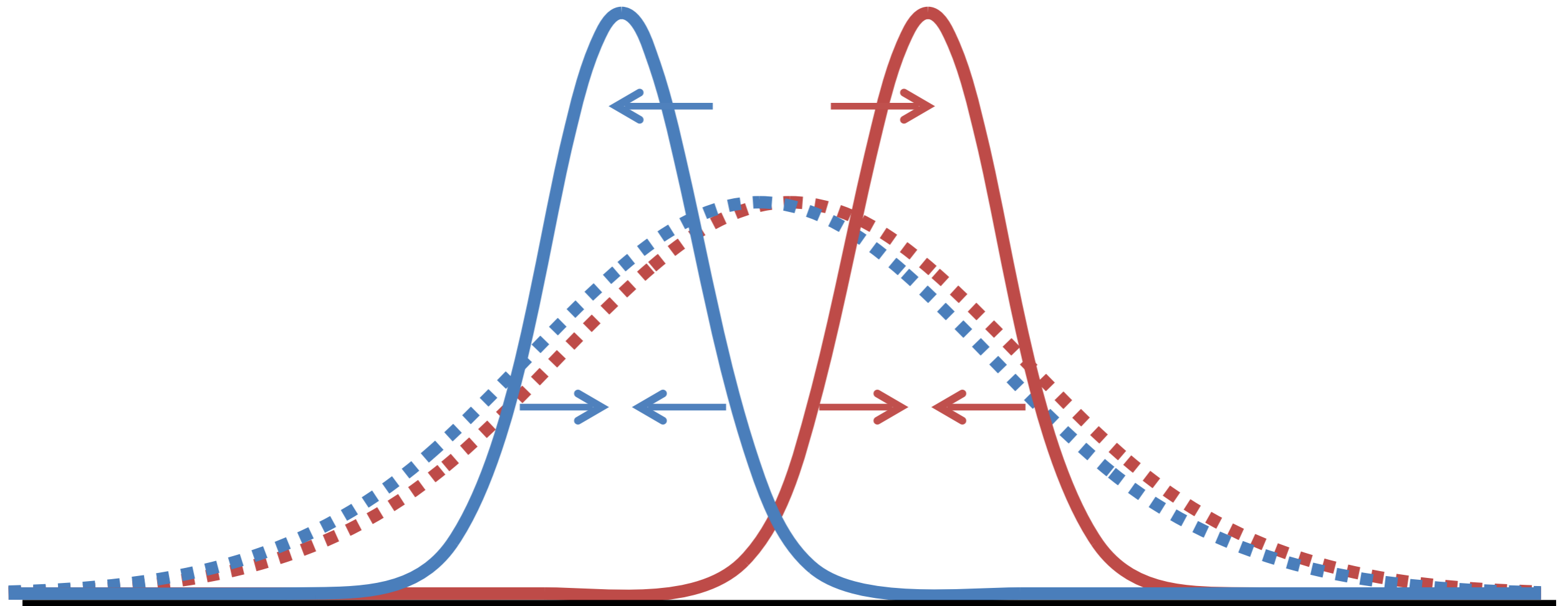


Visualization



Visualization

Observation: **S1** defeats S2



Not pictured: Confidences are separate for each system

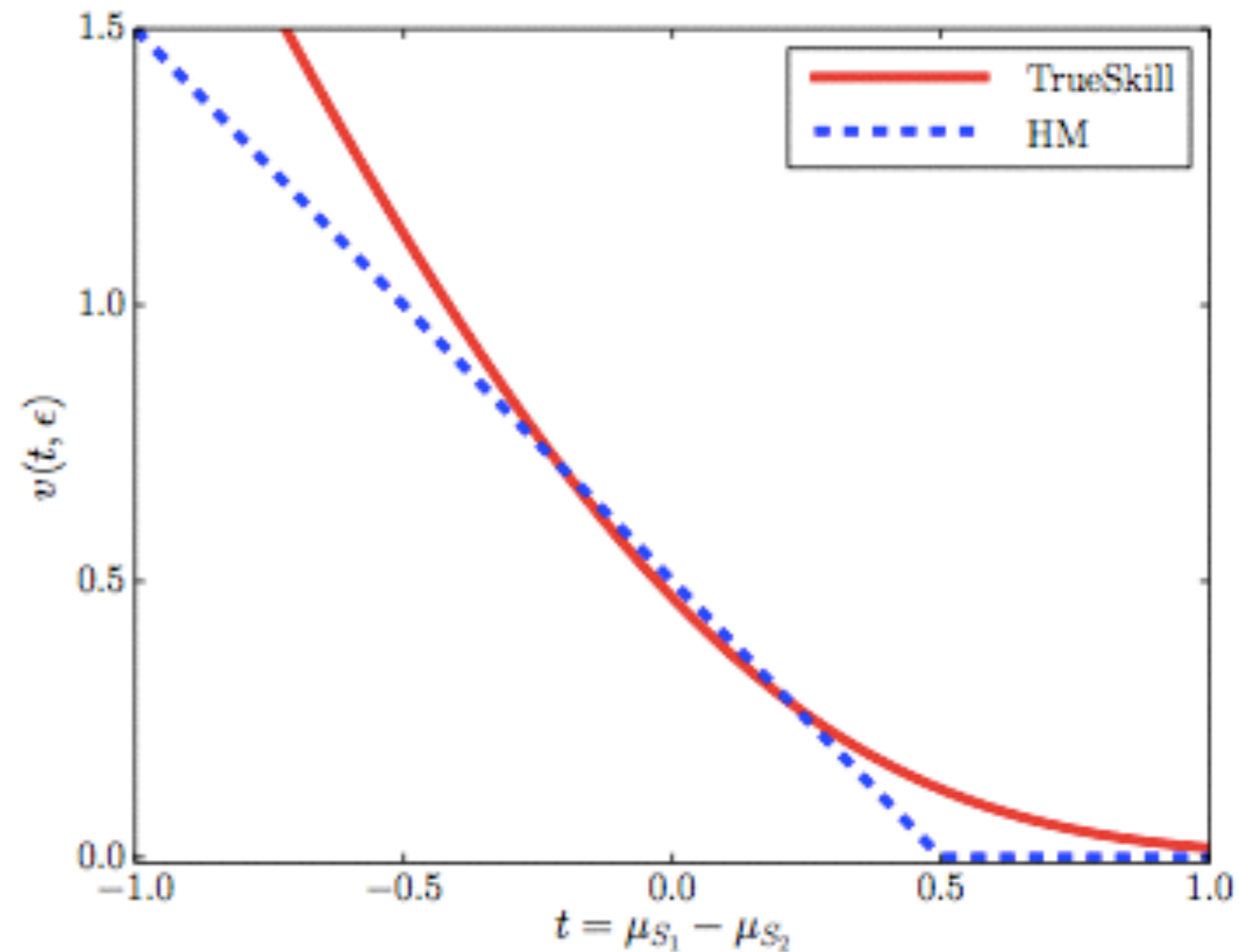
Updating

If S1 defeats S2,

$$\mu_{S_1} = \mu_{S_1} + \frac{\sigma_{S_1}^2}{c} \cdot v\left(\frac{t}{c}, \frac{\epsilon}{c}\right)$$

$$\mu_{S_2} = \mu_{S_2} - \frac{\sigma_{S_2}^2}{c} \cdot v\left(\frac{t}{c}, \frac{\epsilon}{c}\right)$$

outcome surprisal



TrueSkill for MT

- In the MT setting:
 - Each system is a player
 - Each pairwise annotation is a game
- We consider the judgments sequentially, and update the system parameters after each one
- Differences from Xbox:
 - Systems don't improve between games

Procedure

until convergence
 create a new match
 observe the outcome
 update the parameters of both systems

Advantages of TrueSkill

- The system parameter updates reflect how surprising the outcome was
- TrueSkill is an *online* algorithm (as opposed to batch)
 - Instead of sampling system pairs uniformly, we can gather more judgments from systems that are closely matched
 - This presents some potential for reducing the amount of data we need to collect

Partial orderings

- What is the best university in the world?
 - Best is not always well-defined or meaningful
- Instead of total orderings, we use *partial orderings*, which are distinguished by **clusters** of systems that are not distinguished

U.S. News rank	School
#1	Princeton University Princeton, NJ
#2	Harvard University Cambridge, MA
#3	Yale University New Haven, CT
#4	Columbia University New York, NY
#5	Stanford University Stanford, CA
#6	University of California Berkeley Berkeley, CA
#7	University of Michigan Ann Arbor, MI
#8 Tie	University of Pennsylvania Philadelphia, PA
#9	University of Wisconsin-Madison Madison, WI
#10	University of Texas at Austin Austin, TX
#11	University of Washington Seattle, WA
#12	University of Illinois Urbana-Champaign Urbana, IL
#13	University of California San Diego San Diego, CA
#14	University of Wisconsin-Minneapolis Minneapolis, MN
#15	University of Michigan Dearborn Dearborn, MI
#16	University of Wisconsin-Milwaukee Milwaukee, WI
#17	University of Wisconsin-La Crosse La Crosse, WI
#18	University of Wisconsin-River Falls River Falls, WI
#19	University of Wisconsin-Stevens Point Stevens Point, WI
#20	University of Wisconsin-Eau Claire Eau Claire, WI
#21	University of Wisconsin-Oshkosh Oshkosh, WI
#22	University of Wisconsin-Whitewater Whitewater, WI
#23	University of Wisconsin-Superior Superior, WI
#24	University of Wisconsin-Stout Stout, WI
#25	University of Wisconsin-Fox Ochs
#26	University of Wisconsin-Oroquois
#27	University of Wisconsin-Platteville
#28	University of Wisconsin-Stevens Point
#29	University of Wisconsin-Eau Claire
#30	University of Wisconsin-Oshkosh
#31	University of Wisconsin-Whitewater
#32	University of Wisconsin-Superior
#33	University of Wisconsin-Stout
#34	University of Wisconsin-Fox Ochs
#35	University of Wisconsin-Oroquois
#36	University of Wisconsin-Platteville
#37	University of Wisconsin-Stevens Point
#38	University of Wisconsin-Eau Claire
#39	University of Wisconsin-Oshkosh
#40	University of Wisconsin-Whitewater
#41	University of Wisconsin-Superior
#42	University of Wisconsin-Stout
#43	University of Wisconsin-Fox Ochs
#44	University of Wisconsin-Oroquois
#45	University of Wisconsin-Platteville
#46	University of Wisconsin-Stevens Point
#47	University of Wisconsin-Eau Claire
#48	University of Wisconsin-Oshkosh
#49	University of Wisconsin-Whitewater
#50	University of Wisconsin-Superior
#51	University of Wisconsin-Stout
#52	University of Wisconsin-Fox Ochs
#53	University of Wisconsin-Oroquois
#54	University of Wisconsin-Platteville
#55	University of Wisconsin-Stevens Point
#56	University of Wisconsin-Eau Claire
#57	University of Wisconsin-Oshkosh
#58	University of Wisconsin-Whitewater
#59	University of Wisconsin-Superior
#60	University of Wisconsin-Stout
#61	University of Wisconsin-Fox Ochs
#62	University of Wisconsin-Oroquois
#63	University of Wisconsin-Platteville
#64	University of Wisconsin-Stevens Point
#65	University of Wisconsin-Eau Claire
#66	University of Wisconsin-Oshkosh
#67	University of Wisconsin-Whitewater
#68	University of Wisconsin-Superior
#69	University of Wisconsin-Stout
#70	University of Wisconsin-Fox Ochs
#71	University of Wisconsin-Oroquois
#72	University of Wisconsin-Platteville
#73	University of Wisconsin-Stevens Point
#74	University of Wisconsin-Eau Claire
#75	University of Wisconsin-Oshkosh
#76	University of Wisconsin-Whitewater
#77	University of Wisconsin-Superior
#78	University of Wisconsin-Stout
#79	University of Wisconsin-Fox Ochs
#80	University of Wisconsin-Oroquois
#81	University of Wisconsin-Platteville
#82	University of Wisconsin-Stevens Point
#83	University of Wisconsin-Eau Claire
#84	University of Wisconsin-Oshkosh
#85	University of Wisconsin-Whitewater
#86	University of Wisconsin-Superior
#87	University of Wisconsin-Stout
#88	University of Wisconsin-Fox Ochs
#89	University of Wisconsin-Oroquois
#90	University of Wisconsin-Platteville
#91	University of Wisconsin-Stevens Point
#92	University of Wisconsin-Eau Claire
#93	University of Wisconsin-Oshkosh
#94	University of Wisconsin-Whitewater
#95	University of Wisconsin-Superior
#96	University of Wisconsin-Stout
#97	University of Wisconsin-Fox Ochs
#98	University of Wisconsin-Oroquois
#99	University of Wisconsin-Platteville
#100	University of Wisconsin-Stevens Point

#8 Tie	University of Pennsylvania Philadelphia, PA
--------	--

Simulating Human Judgment in Machine Translation Evaluation Campaigns (Koehn, 2012)

Computing clusters

- To compute clusters, we use a statistical technique called *bootstrap resampling*
 - Estimate variance by sampling the sample many times and compute statistics over the samples
- We run each model 1,000 times
 - For each system, extract rank from each fold, throw out outliers
 - Use resulting *rank range* to cluster

X
2
2
2
2
X
2
2
2
2
1
2
2
2
1

Hindi–English (WMT 2014)

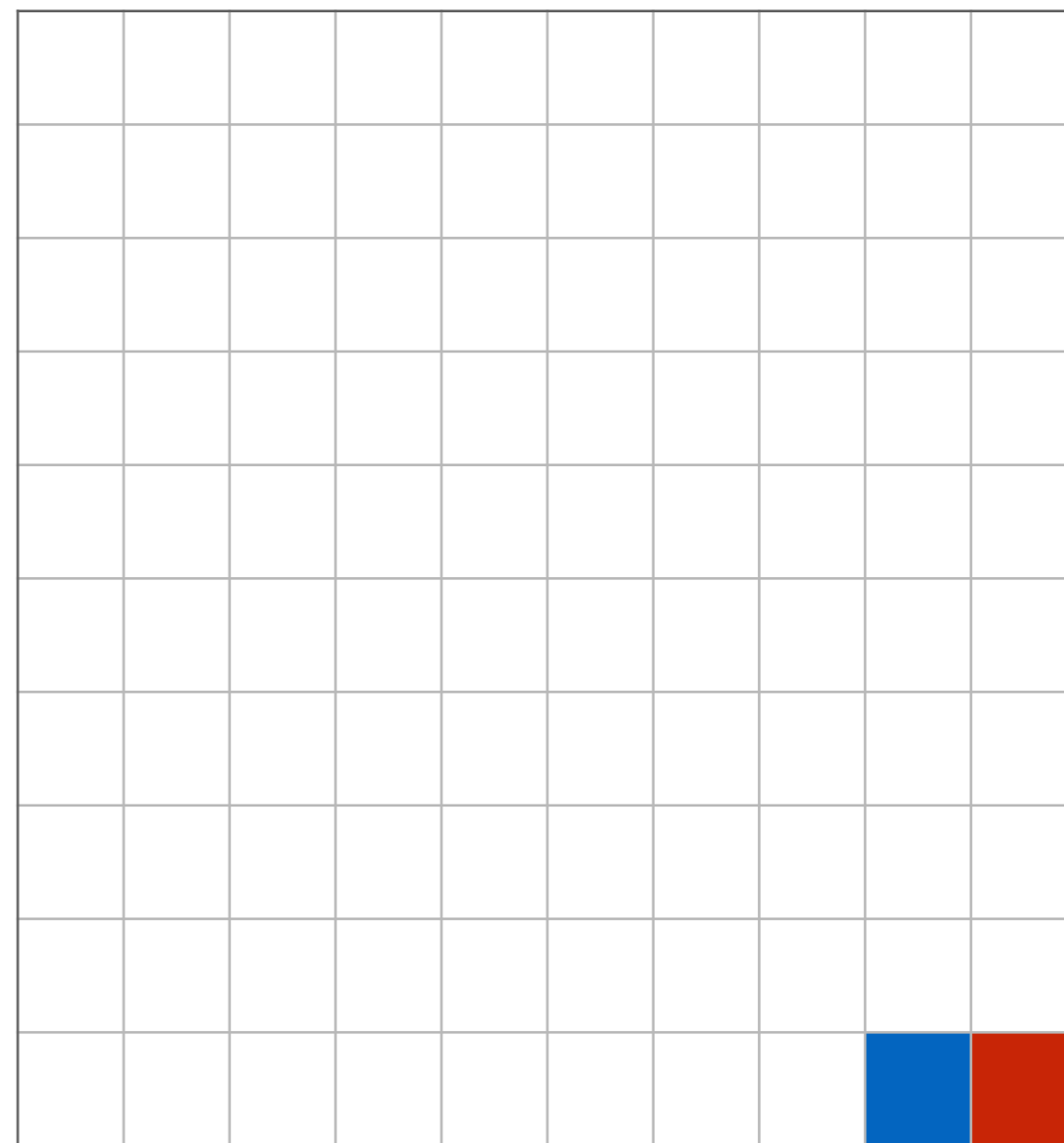
rank range	constrained	unconstrained
1		online-B
2–4	uedin-syntax, cmu	online-A
5	uedin-phrase	
6–7	afri, iit-bombay	
8	dcu-lingo24	
9		iit-hyderabad

Model selection

- We have multiple ways of ranking the systems
 - Expected wins
 - Model of relative ability
 - TrueSkill
- Which is best?
 - Which one does the best job of making predictions?

Model selection

- Experimental setup
 - Split the complete data into 100 folds
 - For each fold
 - Build a model on the other 99 folds
 - Compute accuracy on the current fold
 - Report average accuracy across all folds



*Dataset: 328k judgments
10 language pairs*

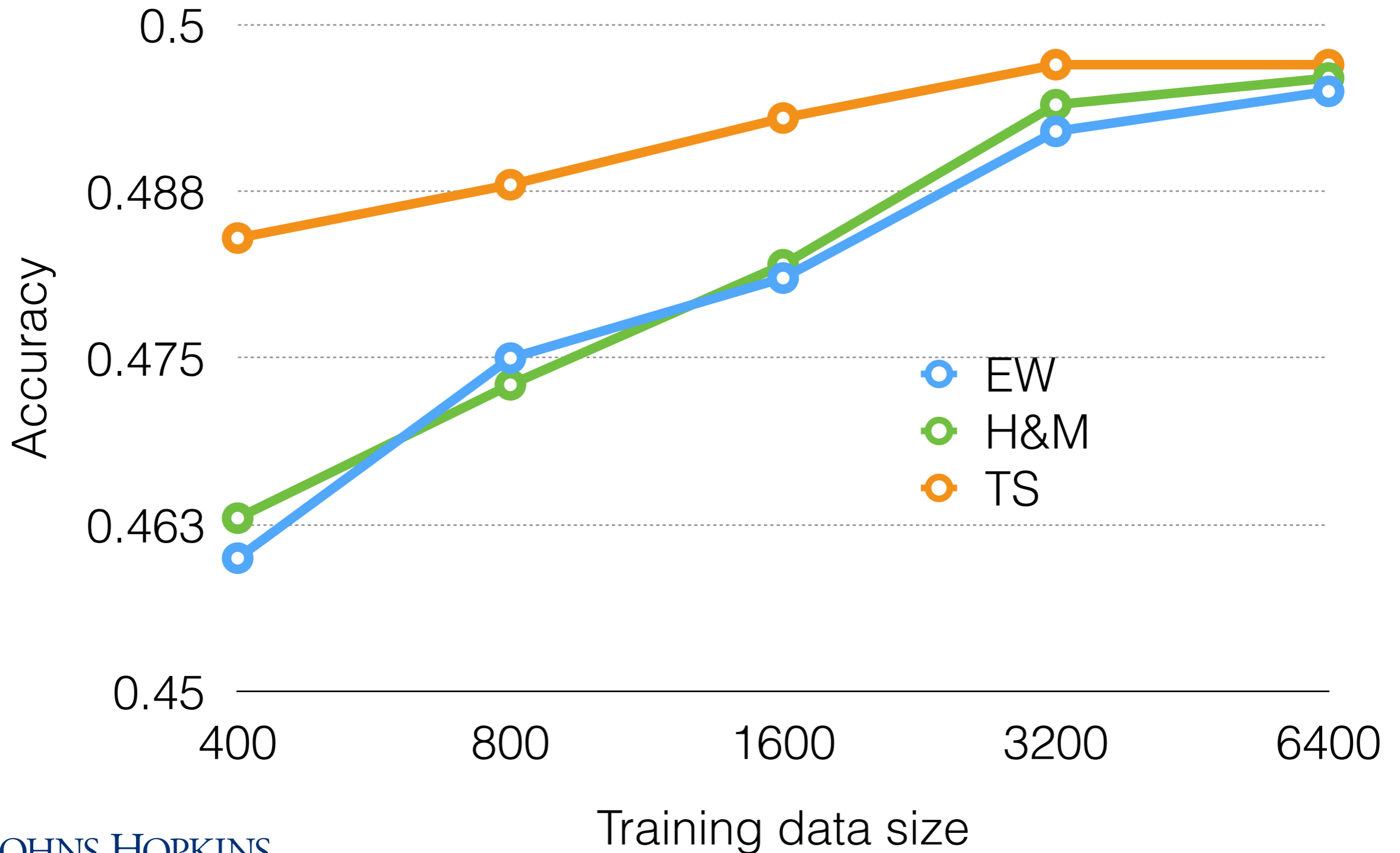
Results

Task	EW	HM	TS	Oracle
Czech–English	40.4	41.1	41.1	41.2
English–Czech	45.3	45.6	45.9	46.8
French–English	49.0	49.4	49.3	50.3
English–French	44.6	44.4	44.7	46.0
German–English	43.5	43.7	43.7	45.2
English–German	47.3	47.4	47.2	48.2
Hindi–English	62.5	62.2	62.5	62.6
English–Hindi	53.3	53.7	53.5	55.7
Russian–English	47.6	47.7	47.7	50.6
English–Russian	46.5	46.1	46.4	48.2
MEAN	48.0	48.1	48.2	49.2

Analysis

- The different methods don't have that much of an effect (surprising?)
 - In fact, the ordering of systems was exactly the same for eight of the language pairs
- However, this hides the amount of data used

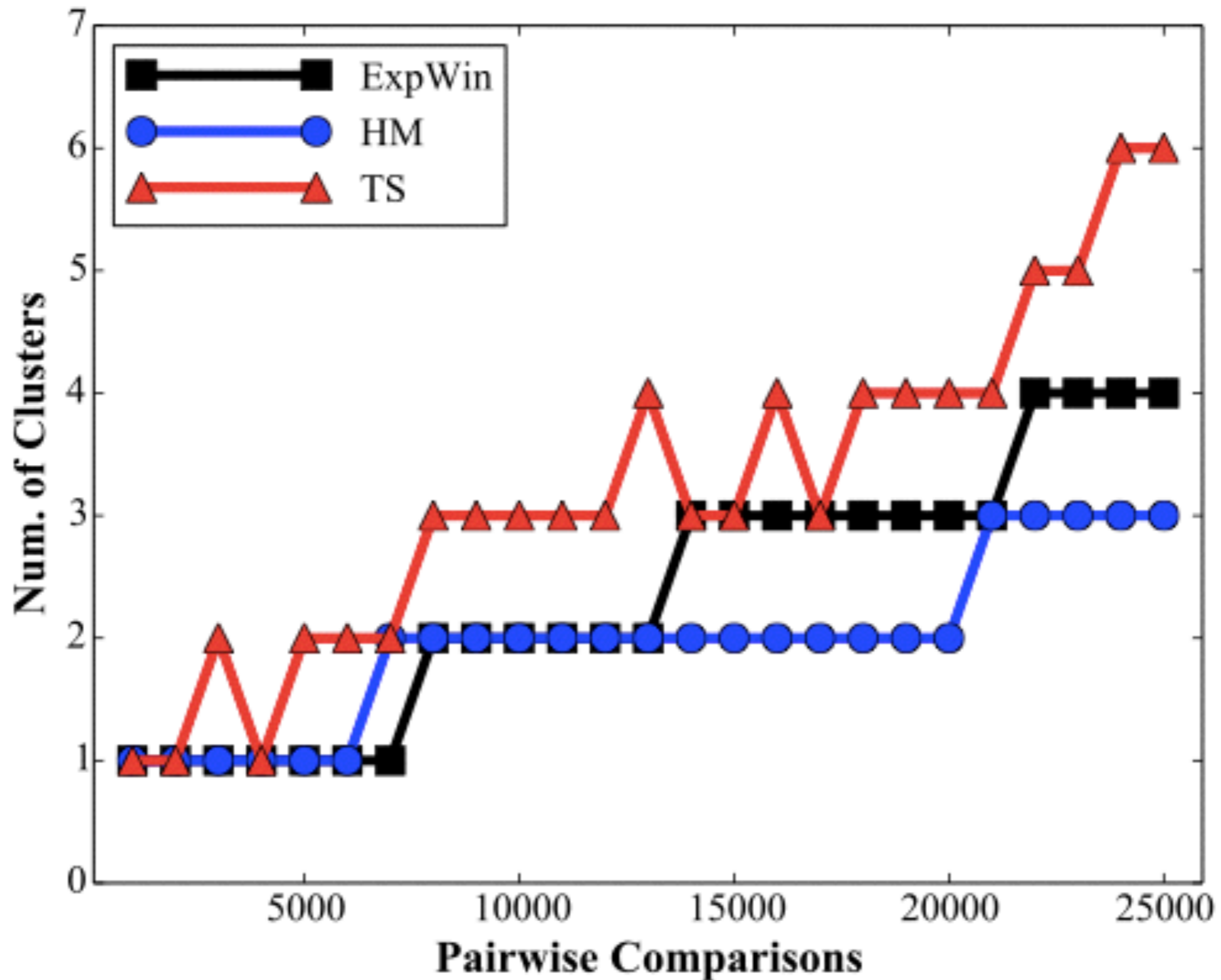
Data requirements



Analysis

- The different methods don't have that much of an effect (surprising?)
 - In fact, the ordering of systems was exactly the same for eight of the language pairs
- However, this hides the amount of data used
 - TrueSkill needs much less data
 - Also has much smaller variance (so we get tighter clusters)

Cluster counts



Summary

- There are many ways of producing the human ranking, from simple models to more elegant ones
- We use the model's ability to predict unseen data as a test of how good it is
 - There are many dimensions to goodness, including accuracy and data requirements
- Translation quality is inherently subjective and task-specific
 - Publishing clusters is a step towards capturing this

