

# Goal: Learn *inflection* $\leftrightarrow$ *root* mappings for world's languages with and without direct supervision

Definition: A **supervised** learning algorithm uses training data from generalized rules can be formed.

Inflection	Root	Part of Speech
<b>English</b>		
swims	swim	3S Present Indicative
swimming	swim	Gerund
swam	swim	1S Past Indicative
swum	swim	Participle
<b>French</b>		
abrège	abrèger	1S Present Indicative
abrègent	abrèger	3P Present Indicative
abrègerai	abrèger	1S Future Indicative
conçu	concevoir	1P Future Anterior Indicative
crois	croire	1S Present Indicative
croyaient	croire	3P Imperfect Indicative
<b>Irish</b>		
thóg	tóg	1S Past Indicative
thógadh	tóg	2P Imperfect Indicative
thógaidís	tóg	3P Imperfect Indicative
adhairim	adhair	1S Present Indicative
d'adhairfeá	adhair	2S Conditional
<b>Dutch</b>		
bood	bieden	3S Past Indicative
gebiedt	bieden	2S Present Indicative
geboden	bieden	1S Past Perfect Indicative
verbod	verbieden	3S Past Indicative
verbiedt	verbieden	2S Present Indicative
aangeboden	aanbieden	1S Past Perfect Indicative

Inflection	Root	Part of Speech
<b>Turkish</b>		
sandınız	sanmak	2P Past Def.Ind.Pos.Int.
sanaydınız	sanmak	2P Past Nar.Sub.Pos.Sta.
sanmayaydınız	sanmak	2P Past Nar.Sub.Neg.Sta.
sanmalıydınız	sanmak	2P Past Nar.Nec.Pos.Sta.
sanmalıymışsınız	sanmak	2P Past Dub.Nec.Pos.Sta.
sanmalıymışsınız	sanmak	2P Past Dub.Nec.Pos.Sta.
sanmamalıymışsınız	sanmak	2P Past Dub.Nec.Neg.Sta.
<b>Tagalog</b>		
gugupitin	gupit	Indicative OF Cont.
pagugupitin	gupit	Causative A <sub>2</sub> F Cont <sub>1</sub>
paggugupitin	gupit	Causative A <sub>2</sub> F Cont <sub>2</sub>
papaggupitin	gupit	Causative A <sub>2</sub> F Inf.
papaggugupitin	gupit	Causative A <sub>2</sub> F Cont <sub>3</sub>
<b>Swahili</b>		
ninaagua	agua	1S Present Indicative
unaagua	agua	2S Present Indicative
niliagua	agua	1S Past Indicative
uliagua	agua	2S Past Indicative
nitaagua	agua	1S Future
utaagua	agua	2S Future
<b>Arabic</b>		
katab	ktb	Active "write"
kattab	ktb	Active "cause to write"
ktutib	ktb	Passive "write"

# Inflectional morphological phenomenon\*

\*from a computational linguistist's P.O.V. & using orthography instead of phonology

<b>affixation</b>	prefixation:	geuza	→	<b>m</b> ligeuza	(Swahili)
	suffixation:	adhair	→	adhair <b>im</b>	(Irish)
	circumfixation:	mischen	→	<b>g</b> emisch <b>t</b>	(German)
	infixation:	palit	→	<b>p</b> umalit	(Tagalog)
<b>point-of-affixation stem changes</b>		placer	→	plaça	(French)
	elision:	close	→	closing	(English)
	gemination:	stir	→	stirred	(English)
	voicing:	zwerft	→	zwerven	(Dutch)
<b>vowel harmony</b>		abartmak	→	abartmasanız	(Turkish)
		addetmek	→	addetmeseniz	(Turkish)
<b>internal vowel shift</b>		afbryde	→	afbrød	(Danish)
		skrike	→	skreik	(Norwegian)
		sleep	→	slept	(English)
<b>agglutination</b>	agglutination:	ev	→	ev <b>de</b>	(Turkish)
	agglutination:		→	evde <b>ki</b>	
	agglutination:		→	evdek <b>iler</b>	
<b>and reduplication</b>	reduplication:	<b>g</b> upit	→	<b>g</b> ugupit	(Tagalog)
	agglutination:		→	<b>ig</b> ugupit	
	agglutination:		→	<b>ip</b> agugupit	
	agglutination:		→	<b>ipin</b> agugupit	
<b>root and pattern</b>	reduplication:	rumah	→	<b>rumah</b> rumah	(Malay)
	reduplication:	ibu	→	<b>ibu</b> ibu	
<b>highly irregular forms</b>		ktb	→	kateb	(Arabic)
		ktb	→	kattab	
		fi	→	erai	(Romanian)
	jānā	→	gayā	(Hindi)	
	eiga	→	áttum	(Icelandic)	

---

## Task Definition

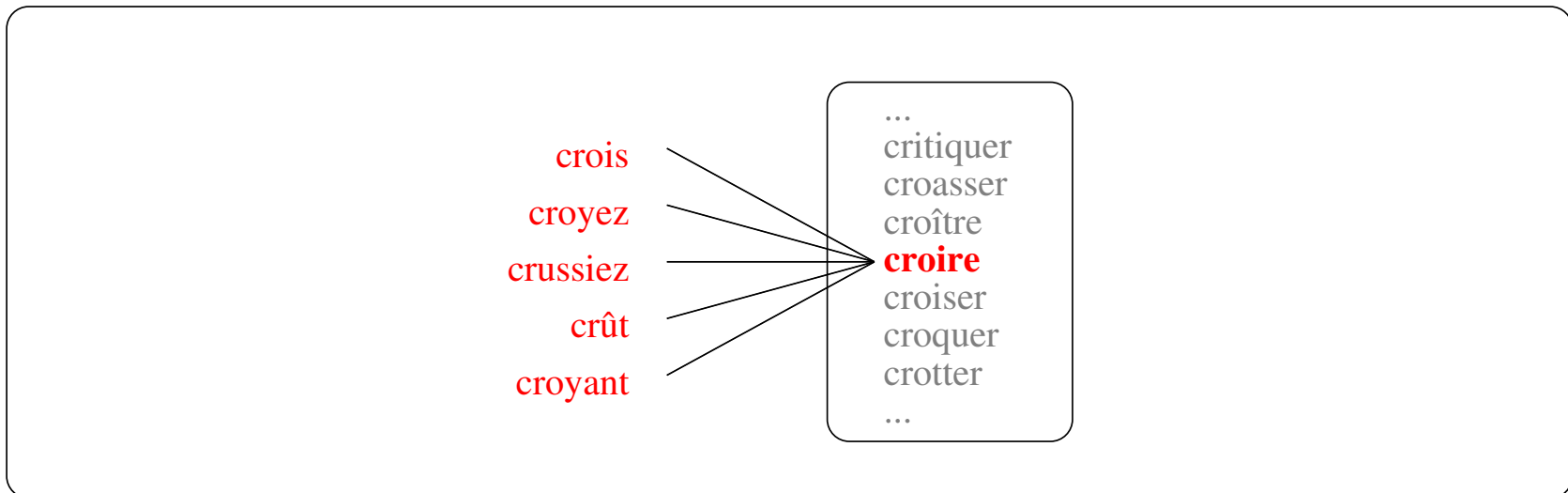
	Morphological Analysis	Morphological Generation
Input	<u>inflection</u> <i>burned</i>	<u>root, part of speech</u> <i>burn, VBD</i>
Output	<u>root, (optional) part of speech</u> <i>burn, VBD [Past Indicative]</i> <i>burn, VBN [Past Participle]</i>	<u>inflection</u> <i>burnt</i> <i>burned</i>

- Notice that both morphological analysis and morphological generation can often generate multiple correct answers.
- When the part-of-speech is omitted in morphological analysis, the task is often referred to as morphological stemming.

---

# Major applications of computational morphology

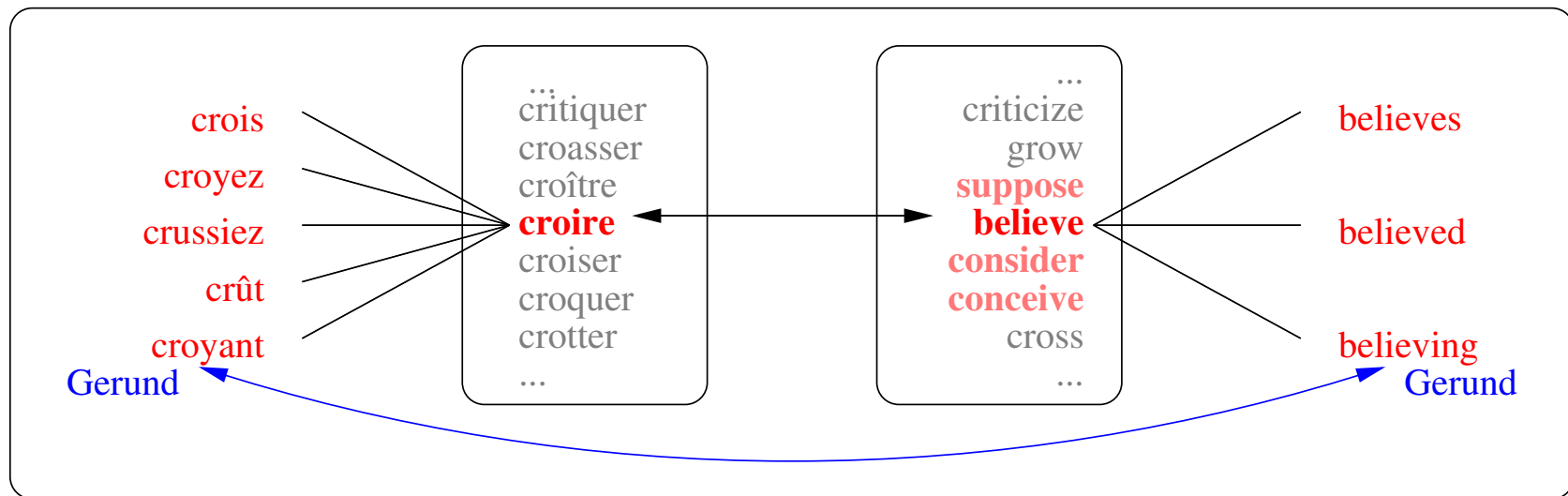
- Information retrieval
  - dimensionality reduction



- Machine translation
  - translation lexicon access, dimensionality reduction for contextual features, fine grained part of speech
- Other applications: parsing, word sense disambiguation, text generation, part of speech tagging

# Major applications of computational morphology

- Information retrieval
  - dimensionality reduction
- Machine translation
  - translation lexicon access, dimensionality reduction for contextual features, fine grained part of speech



- Other applications: parsing, word sense disambiguation, text generation, part of speech tagging

---

## Alignment Paradigm

- Prior approaches have focused (almost) exclusively on learning string transductions.
- This makes learning irregular morphology, and pairs such as the following, difficult:

sang	↔	sing
singed	↔	singe

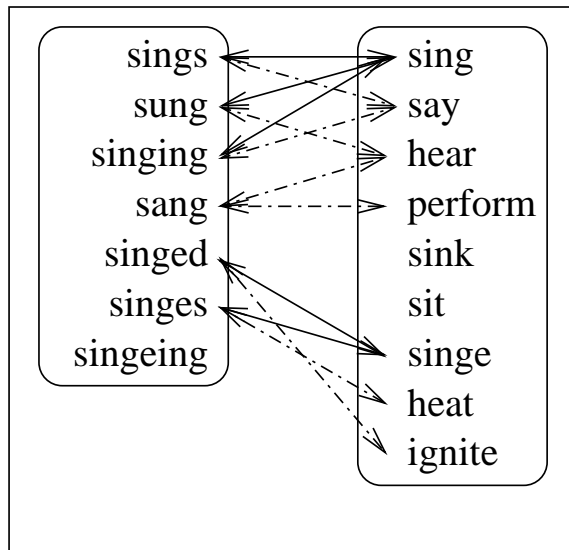
- How can we learn that the past tense of *sing* isn't *singed*?
- Possible answer: a large amount of information for inflection-root mappings is available outside the string:

Context similarity		Distributional similarity	
sing	songs	sing	1204
sang	songs	sang	1427
singe	hair	singe	2
singes	hair	singes	9

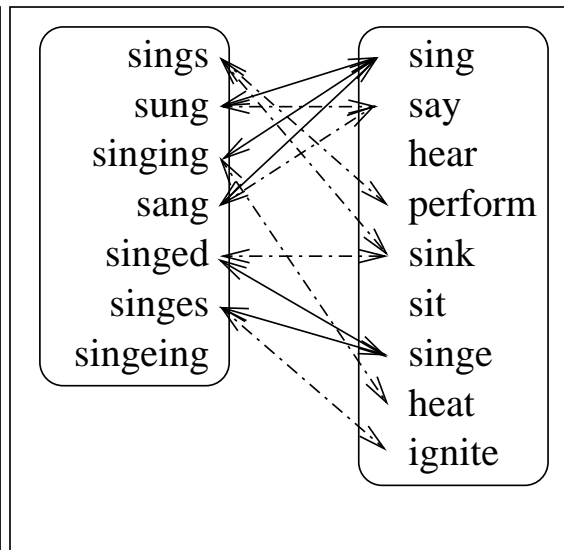
## Alignment Paradigm

- Treat morphological analysis as an inflection-root mapping problem.
- Use multiple similarity measures and find a consensus answer.
  - Positionally weighted contextual similarity
  - Distributional similarity (frequency)
  - Weighted Levenshtein similarity (string-edit distance)
  - Bilingual projection

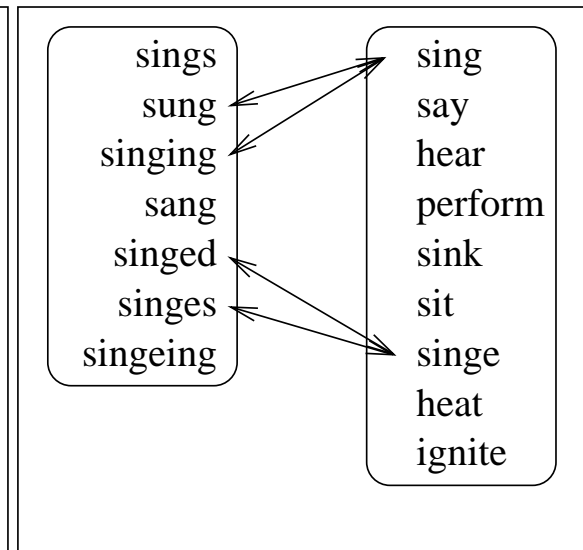
Contextual Similarity



Distributional Similarity



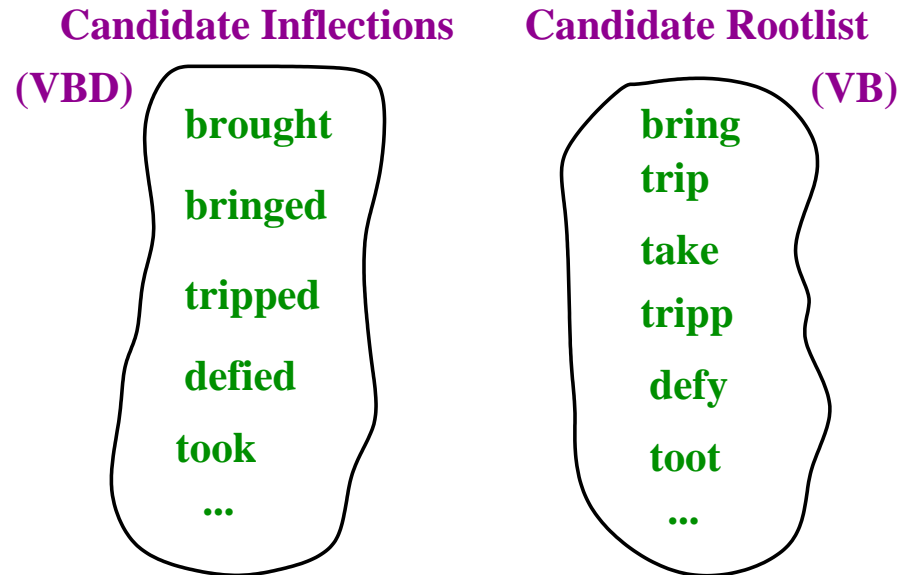
Combined Similarity



---

## Resource Assumptions

- Noisy wordlists of inflection and root candidates for the language:



- Canonical suffixes of the language (optional):

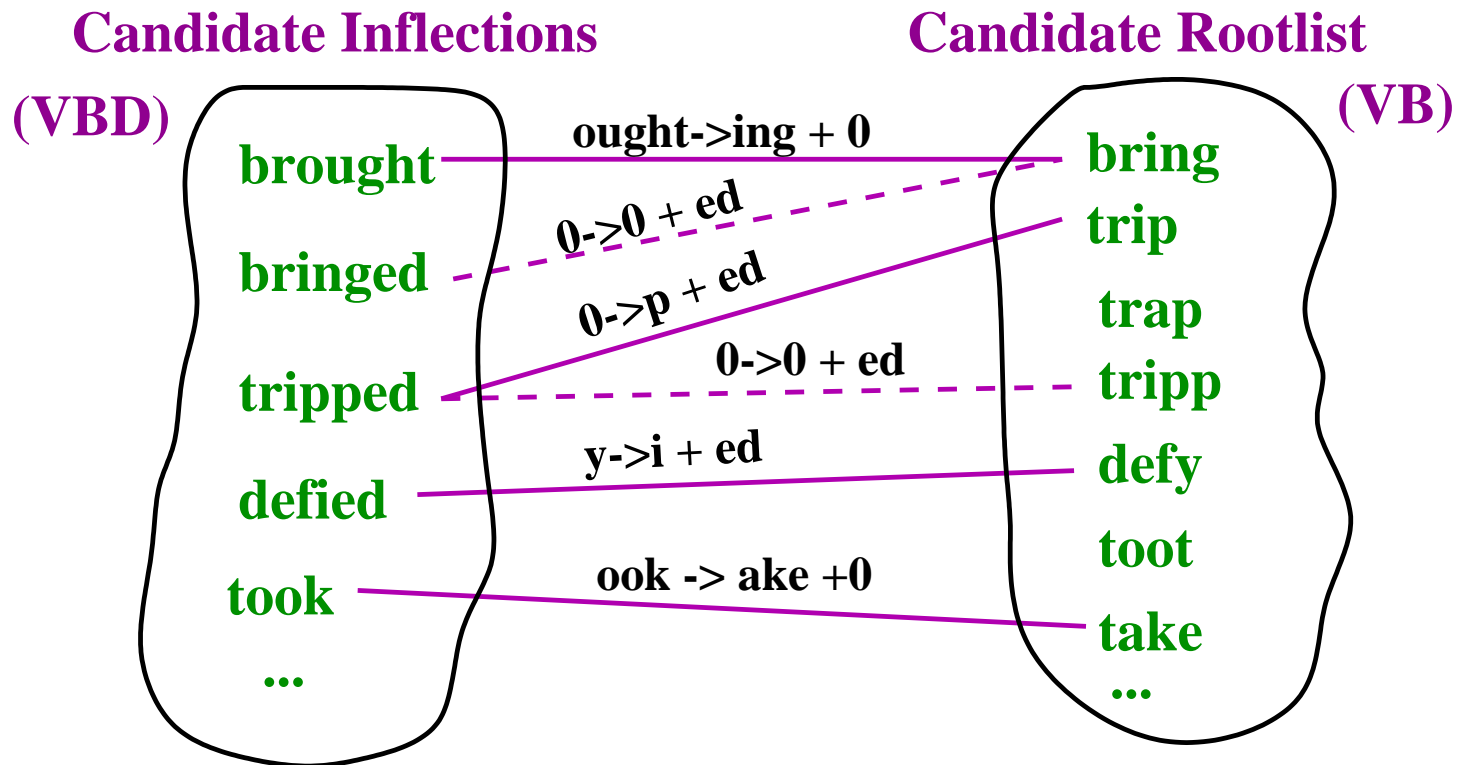
Part of Speech	VB	VBD	VBZ	VBG	VBN
Canonical Suffixes	+ $\epsilon$	+ed + $\epsilon$	+s	+ing	+en +ed + $\epsilon$

- Lots of plain text



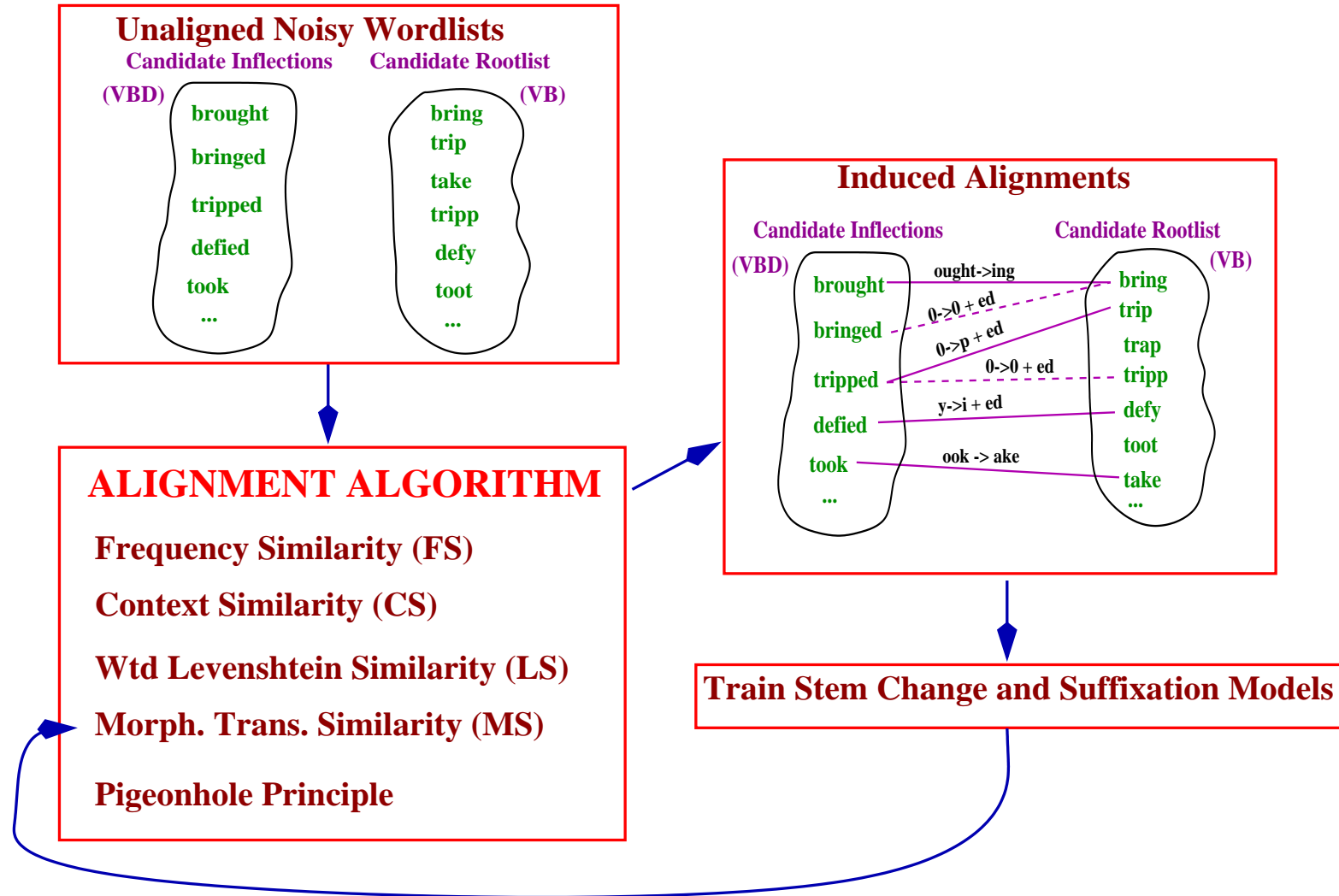
# Approach

- Treat analysis as probabilistic alignment over large wordlists.



- Use these alignments to train string transduction models

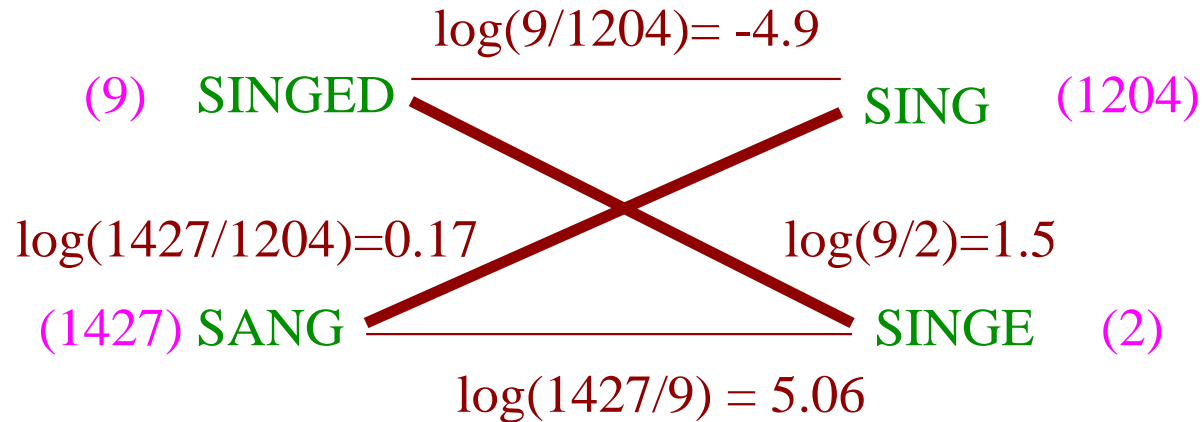
# Iterative Bootstrapping of Similarity Models



---

## Distributional Similarity Models

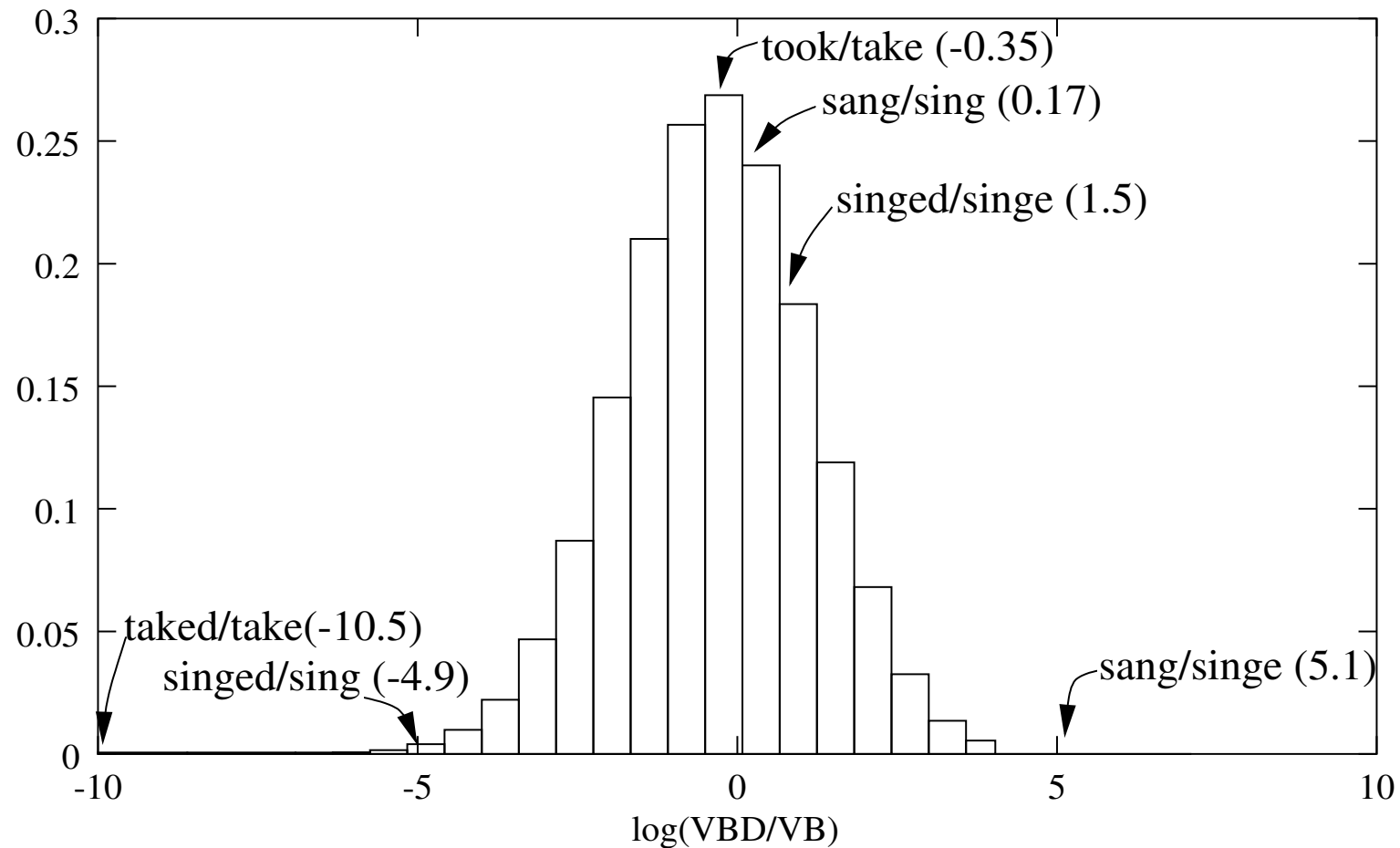
Favor inflection/root alignments with “good match” of frequencies



- How to quantify “good match”?
- How to penalize divergence?

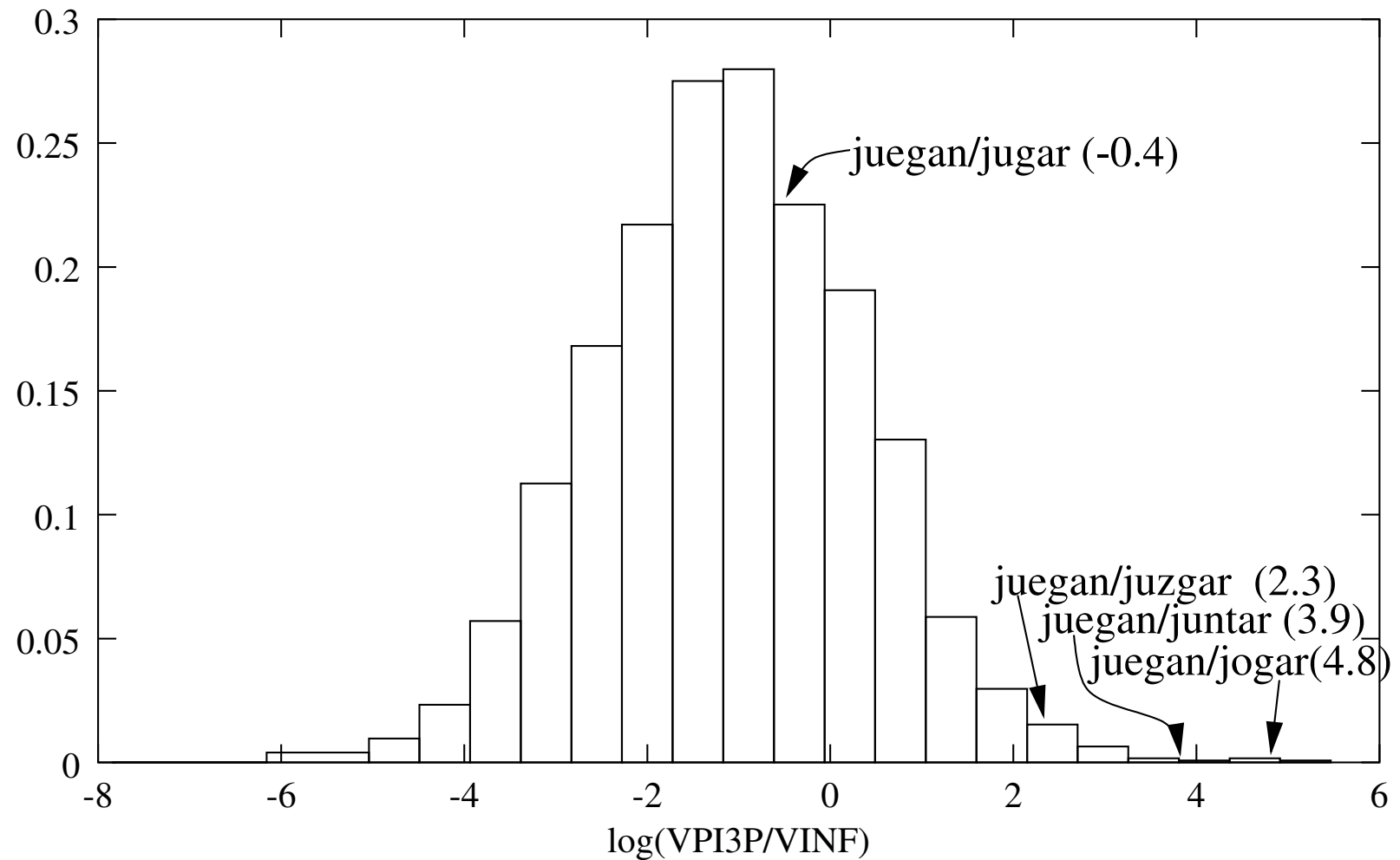
# Distributional Similarity Models

Use empirical distributions of frequency ratios to measure “goodness” of fit.



$$E[\log(\text{VBD}/\text{VB})] = -0.24$$

## Distributional Similarity Models (Spanish)

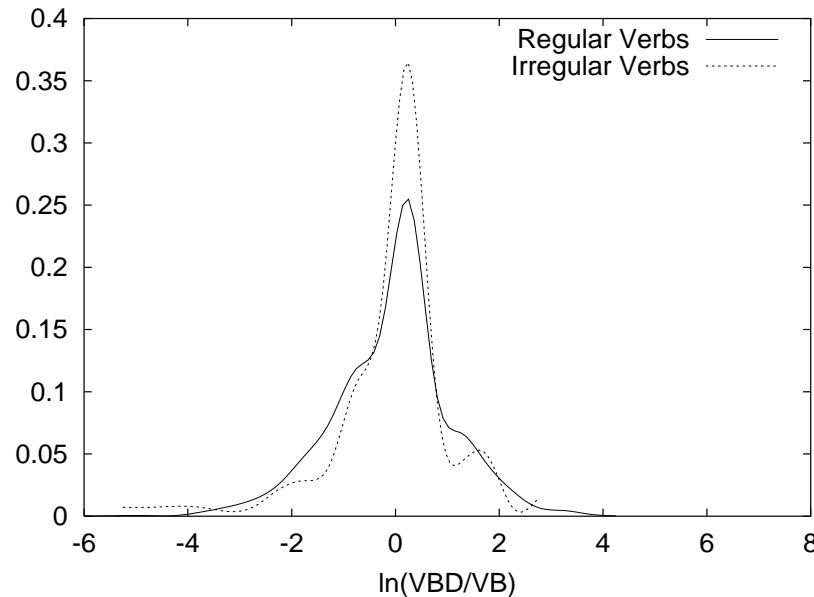


$E[\log(\text{VP13P}/\text{VINF})] = -1.21$

---

# Approximating Full Empirical Ratio Distributions

VerbType	$\frac{VBD}{VB}$	$\frac{VBG}{VB}$	Avg. Lemma Freq
Regular	.847	.746	861
Irregular	.842	.761	17406



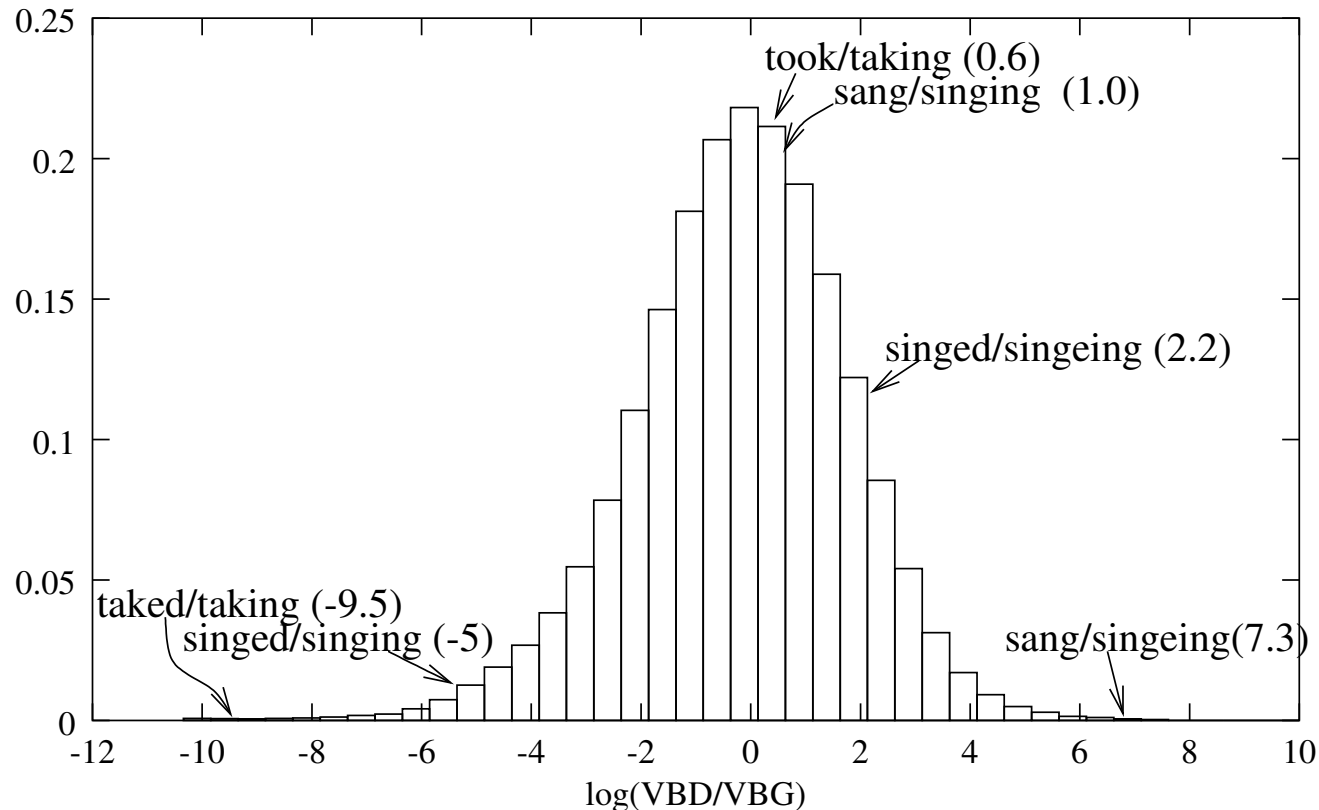
- Tense distribution of verb inflections not correlated with their regularity or stem changing properties
- Estimate initial empirical distributions from most confidently alignable pairs (often regular inflections) from other models

---

## Multiple Ratio Estimators of Robustness

Comparing the distributional frequency directly with other inflections are also informative

e.g.  $f(\text{VBD})/f(\text{VBG})$  or  $f(\text{VBD})/f(\text{VBZ})$

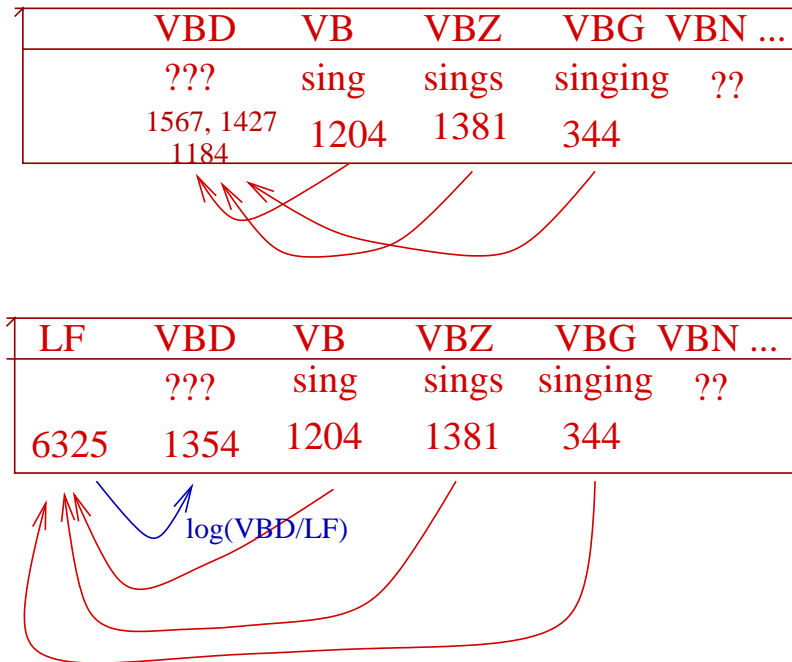


- Any one individual inflection ratio can vary widely due to poor alignment or idiosyncratic tense usage

---

## Reducing Frequency Model Dimensionality

- Use hidden variable of Estimated Lemma Frequency (LF) to capture and smooth multiple pairwise indicators



- Decreases pairwise model dimensionality ( $T^2 \rightarrow T$ )
- Improves robustness and coverage (all important for highly inflected languages)



---

## Context Similarity

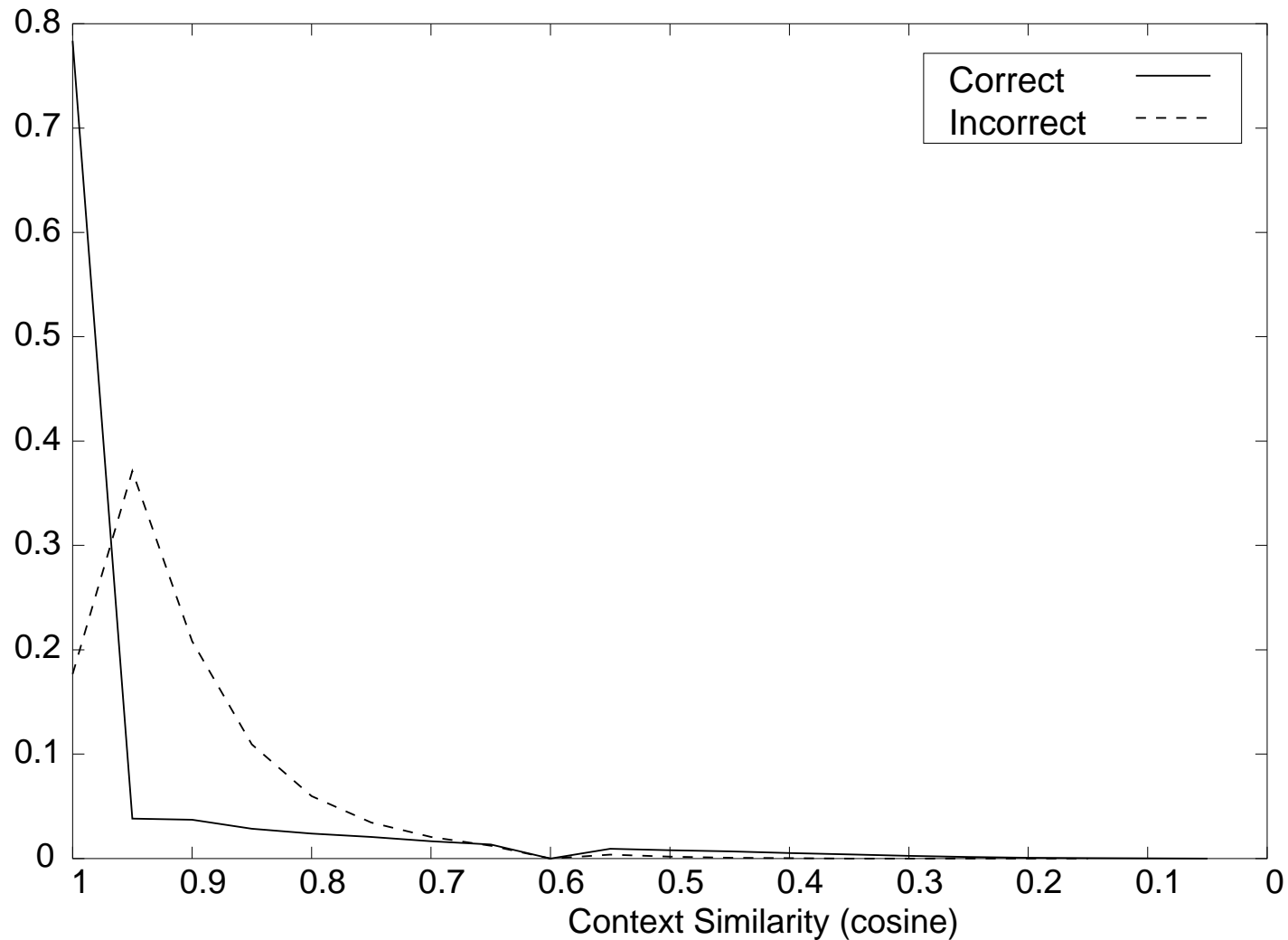
Measure cosine similarity between aggregate, position-weighted context vectors

	hands	head	baby	violently	himself	deer	away
shook	128	103	21	17	-	-	1
shake	151	98	8	12	-	-	-
shoot	-	-	-	-	56	8	1
shoo	-	-	-	-	-	-	6

- Pool of basic regular expressions to locate potential salient positions.
- Regex choice and relative weighting optimized empirically on strongest alignments from other models.

---

## Context similarity distributions for correctly and incorrectly aligned inflection-root pairs (French)



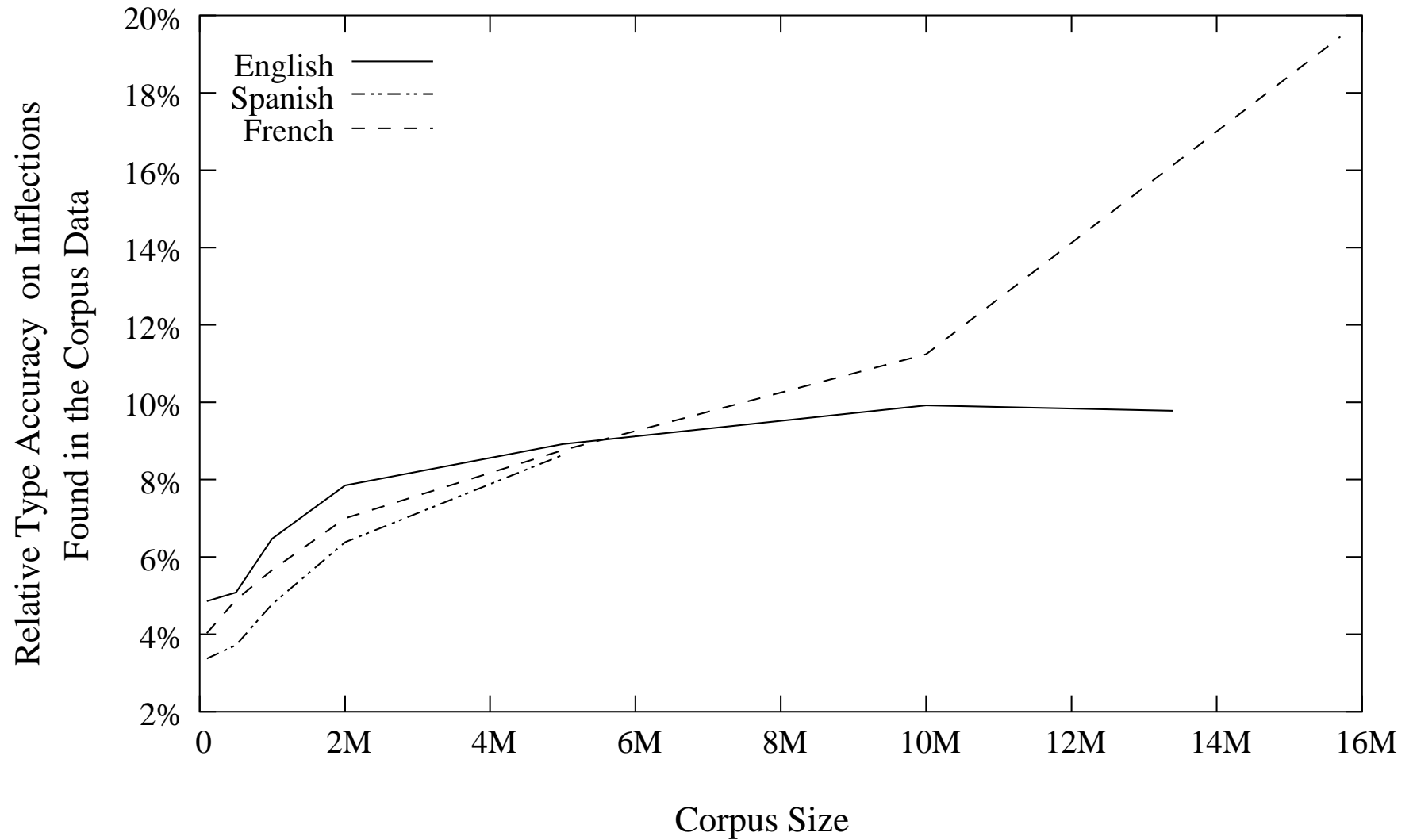
## Context model sensitivity to window position

Language	Left 6...0	5...1	4...2	Center 3...3	2...4	1...5	Right 0...6
<b>S-V-O</b>							
Portuguese	12.05%	21.38%	27.55%	26.58%	29.01%	29.01%	<b>32.87%</b>
Estonian	31.87%	41.43%	44.01%	42.20%	44.28%	<b>44.70%</b>	32.21%
<b>Free / S-V-O</b>							
Russian	19.91%	41.08%	47.35%	40.91%	<b>47.90%</b>	47.39%	47.35%
<b>Verb Second (V2)</b>							
German	9.97%	11.34%	13.09%	<b>14.78%</b>	13.78%	13.34%	9.96%
<b>S-O-V</b>							
Turkish	<b>52.66%</b>	49.06%	48.15%	44.40%	47.03%	45.41%	25.28%
Basque	<b>25.88%</b>	21.35%	21.91%	19.75%	21.66%	19.81%	6.44%

- Bag of 6 words surrounding target word
  - 6...0: word<sub>1</sub> word<sub>2</sub> word<sub>3</sub> word<sub>4</sub> word<sub>5</sub> word<sub>6</sub> **target**
  - 3...3: word<sub>1</sub> word<sub>2</sub> word<sub>3</sub> **target** word<sub>4</sub> word<sub>5</sub> word<sub>6</sub>
  - 0...6: **target** word<sub>1</sub> word<sub>2</sub> word<sub>3</sub> word<sub>4</sub> word<sub>5</sub> word<sub>6</sub>

Language	Left 6...0	Center 3...3	Right 0...6
<b>S-V-O</b>			
Spanish	6.37%	20.31%	<b>29.38%</b>
Portuguese	12.05%	26.58%	<b>32.87%</b>
French	9.08%	38.40%	<b>45.60%</b>
Italian	3.69%	9.99%	<b>14.98%</b>
Romanian	10.42%	18.71%	<b>20.86%</b>
English	13.25%	21.98%	<b>25.67%</b>
Danish	7.21%	24.61%	<b>34.59%</b>
Swedish	2.09%	10.36%	<b>18.69%</b>
Icelandic	10.93%	23.43%	<b>29.98%</b>
Estonian	31.87%	<b>42.20%</b>	32.21%
Finnish	5.40%	12.09%	<b>12.15%</b>
Tagalog	10.10%	15.08%	<b>17.08%</b>
Swahili	8.63%	8.68%	<b>11.02%</b>
<b>Free / S-V-O</b>			
Czech	3.30%	<b>11.05%</b>	11.02%
Polish	8.16%	18.91%	<b>20.87%</b>
Russian	19.91%	40.91%	<b>47.35%</b>
<b>Verb Second (V2)</b>			
German	9.97%	<b>14.78%</b>	9.96%
Dutch	11.78%	<b>16.32%</b>	15.50%
<b>S-O-V</b>			
Turkish	<b>52.66%</b>	44.40%	25.28%
Basque	<b>25.88%</b>	19.75%	6.44%

## Context model sensitivity to corpus size



---

## Weighted Levenshtein Similarity

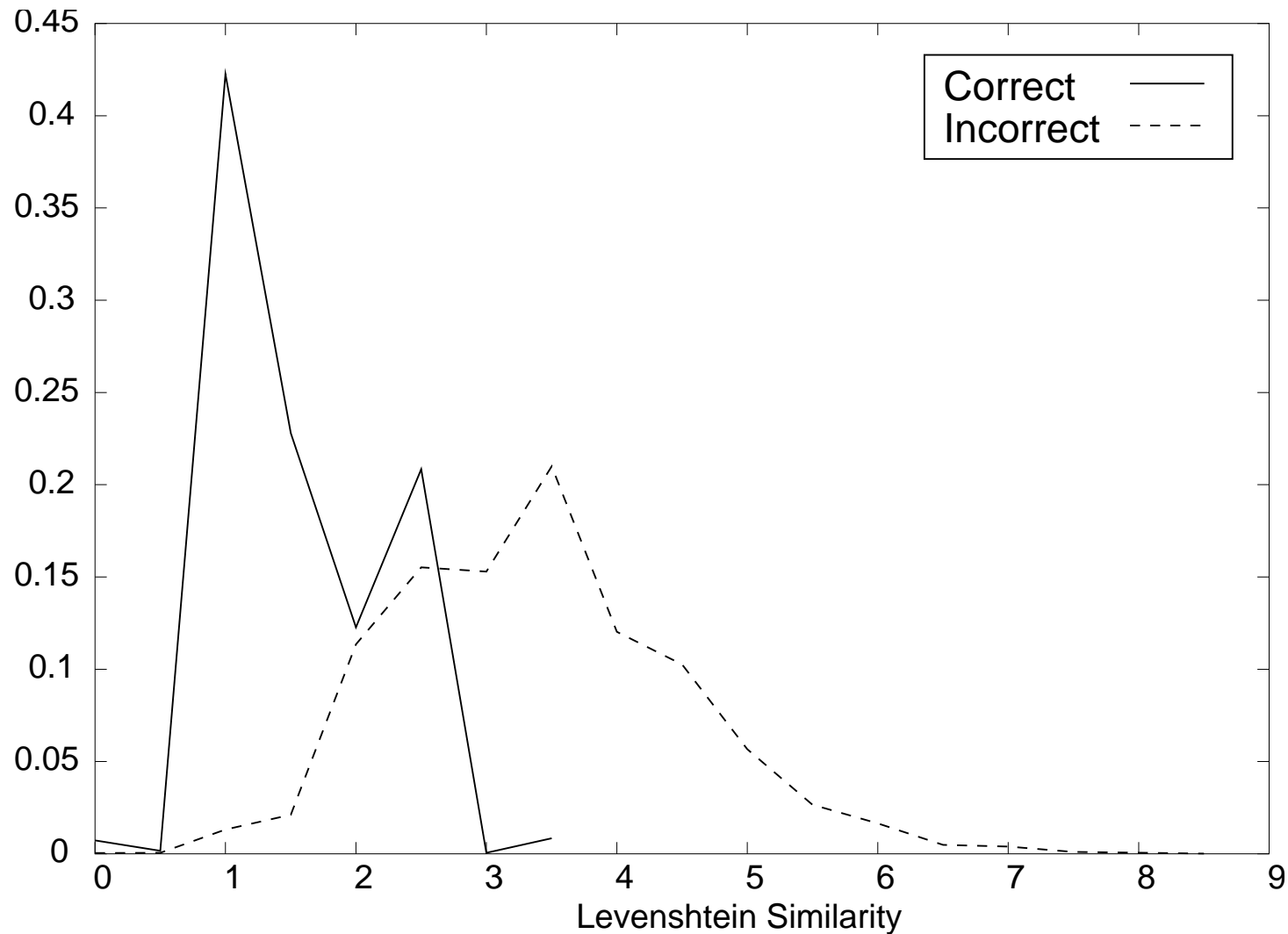
- Measure of string edit distance
- Transition cost matrix weighted by relative rarity of letter change in current paired data

	a	o	ue	m	n	...
a	0	$\delta_1$	$\delta_2$	$\delta_4$	$\delta_4$	...
o	$\delta_1$	0	$\delta_2$	$\delta_4$	$\delta_4$	...
ue	$\delta_2$	$\delta_2$	0	$\delta_4$	$\delta_4$	...
m	$\delta_4$	$\delta_4$	$\delta_4$	0	$\delta_3$	...
n	$\delta_4$	$\delta_4$	$\delta_4$	$\delta_3$	0	...
...	...	...	...	...	...	...

- Initially set to 4 basic parameters for V-V, V<sup>+</sup>-V<sup>+</sup>, C-C, C-V<sup>+</sup> in current paired data
- Cost matrix re-estimated on subsequent alignments

---

## Levenshtein similarity distributions for correctly and incorrectly aligned inflection-root pairs (English)



Language	Prefix Penalty			No Penalty	Suffix Penalty		
	1	0.5	0.25		0.25	0.5	1.0
Spanish	<b>94.13%</b>	93.88%	93.36%	91.62%	63.38%	69.64%	76.08%
Portuguese	<b>96.26%</b>	96.01%	95.38%	93.74%	71.79%	76.24%	80.89%
Catalan	<b>92.01%</b>	91.37%	90.51%	88.02%	71.84%	74.96%	79.17%
Occitan	<b>92.64%</b>	92.40%	91.98%	88.63%	67.28%	71.21%	75.14%
French	<b>93.44%</b>	92.80%	91.78%	88.09%	69.17%	72.05%	76.45%
Italian	<b>95.26%</b>	94.86%	94.37%	91.79%	71.58%	75.70%	80.74%
Romanian	<b>91.14%</b>	90.48%	89.34%	82.67%	46.08%	51.93%	58.64%
Latin	<b>85.35%</b>	84.67%	82.89%	70.60%	21.66%	26.47%	34.02%
English	<b>93.36%</b>	92.38%	89.82%	84.80%	46.96%	54.28%	61.76%
Danish	<b>94.77%</b>	93.31%	92.69%	90.94%	70.44%	76.80%	81.00%
Norwegian	<b>94.47%</b>	93.81%	93.14%	91.40%	72.42%	78.51%	82.96%
Swedish	<b>90.24%</b>	88.93%	87.15%	82.74%	42.88%	51.00%	59.83%
Icelandic	<b>91.33%</b>	90.95%	90.30%	88.65%	62.38%	67.55%	74.11%
Hindi	<b>96.88%</b>	96.48%	<b>96.88%</b>	96.48%	87.50%	87.50%	88.67%
Sanskrit	<b>78.34%</b>	77.39%	75.75%	67.43%	29.68%	34.41%	41.58%
Estonian	<b>81.86%</b>	81.20%	80.70%	78.52%	62.59%	65.71%	69.17%
Tamil	<b>89.61%</b>	87.60%	87.27%	83.42%	62.31%	69.01%	74.54%
Finnish	<b>74.86%</b>	73.57%	71.88%	62.88%	27.35%	32.43%	39.52%
Turkish	<b>95.03%</b>	94.69%	94.09%	89.63%	48.85%	55.83%	64.61%
Uzbek	<b>84.67%</b>	84.43%	83.89%	81.02%	51.19%	55.21%	60.40%
Basque	<b>80.74%</b>	79.85%	78.69%	73.91%	38.45%	44.21%	49.47%
Czech	76.49%	76.40%	76.42%	<b>78.26%</b>	67.13%	70.13%	72.79%
Polish	<b>93.22%</b>	93.02%	92.78%	91.71%	68.54%	73.57%	78.83%
Russian	<b>84.73%</b>	83.41%	82.12%	80.87%	66.33%	69.59%	72.96%
German	<b>91.58%</b>	91.53%	91.39%	91.44%	82.02%	84.04%	86.34%
Dutch	<b>80.49%</b>	80.11%	80.22%	78.08%	69.56%	71.59%	73.79%
Irish	<b>92.89%</b>	<b>92.89%</b>	92.21%	87.75%	48.71%	54.16%	62.18%
Welsh	<b>85.99%</b>	84.81%	83.52%	76.50%	35.41%	42.08%	50.69%
Tagalog	20.27%	23.80%	28.35%	61.03%	72.24%	<b>72.32%</b>	71.40%
Swahili	29.83%	34.94%	41.69%	68.38%	<b>79.95%</b>	79.35%	78.31%
Klingon	24.84%	27.34%	30.46%	98.74%	<b>99.55%</b>	99.17%	99.08%



---

## Using Alignment Models to Bootstrap String Transduction Models

- Combine alignment models to get training data
- Models each have different dynamic ranges
  - Levenshtein Similarity:  $[0, \text{inf}) \Rightarrow$  lower score is better
  - Context Similarity:  $[0, 1] \Rightarrow$  higher score is better
  - Frequency Similarity  $(-\text{inf}, \text{inf})$
- Bidirectional averaged relative rankings

$$\begin{aligned} \text{sim}(\text{root}_i, \text{inf}_j) = & \\ & \lambda_F(\text{rank}_{f_s}(\text{root}_i | \text{inf}_j) + \text{rank}_{f_s}(\text{inf}_j | \text{root}_i)) + \\ & \lambda_C(\text{rank}_{c_s}(\text{root}_i | \text{inf}_j) + \text{rank}_{c_s}(\text{inf}_j | \text{root}_i)) + \\ & \lambda_L(\text{rank}_{l_s}(\text{root}_i | \text{inf}_j) + \text{rank}_{l_s}(\text{inf}_j | \text{root}_i)) + \\ & \lambda_M(\text{rank}_{m_s}(\text{root}_i | \text{inf}_j) + \text{rank}_{m_s}(\text{inf}_j | \text{root}_i)) \end{aligned}$$

- Favors mutual affinity

---

## Choosing $\lambda$ 's for Model Combination

$$\begin{aligned} \text{sim}(\text{root}_i, \text{inf}_j) = & \\ & \lambda_F(\text{rank}_{f_s}(\text{root}_i|\text{inf}_j) + \text{rank}_{f_s}(\text{inf}_j|\text{root}_i)) + \\ & \lambda_C(\text{rank}_{c_s}(\text{root}_i|\text{inf}_j) + \text{rank}_{c_s}(\text{inf}_j|\text{root}_i)) + \\ & \lambda_L(\text{rank}_{l_s}(\text{root}_i|\text{inf}_j) + \text{rank}_{l_s}(\text{inf}_j|\text{root}_i)) + \\ & \lambda_M(\text{rank}_{m_s}(\text{root}_i|\text{inf}_j) + \text{rank}_{m_s}(\text{inf}_j|\text{root}_i)) \end{aligned}$$

- Initially, Levenshtein model will produce the cleanest training data for bootstrapping  $\Rightarrow$  set  $\lambda_L \gg \lambda_C$  and  $\lambda_F$
- The output of the string transduction model can then be used to refine the alignment models
- As context and frequency models are refined,  $\lambda_C$  and  $\lambda_F$  are increased relative to  $\lambda_L$

## Comparison of Models (English)

Combination of Similarity Models	# of Iterations	All Words (3888)	Highly Irregular (128)	Regular (1877)	Semi-Regular (1883)
FS ( <i>Frequency Sim</i> )	(Iter 1)	9.8	18.6	8.8	10.1
LS ( <i>Levenshtein Sim</i> )	(Iter 1)	31.3	19.6	20.0	34.4
CS ( <i>Context Sim</i> )	(Iter 1)	28.0	32.8	30.0	25.8
CS+FS	(Iter 1)	32.5	64.8	32.0	30.7
CS+FS+LS	(Iter 1)	71.6	76.5	71.1	71.9
CS+FS+LS+MS	(Iter 1)	96.5	74.0	97.3	97.4
<b>CS+FS+LS+MS</b>	(Conv)	<b>99.2</b>	<b>80.4</b>	<b>99.9</b>	<b>99.7</b>
<b>Mooney&amp;Califf</b>		82.5	5.0	100.0	84.0

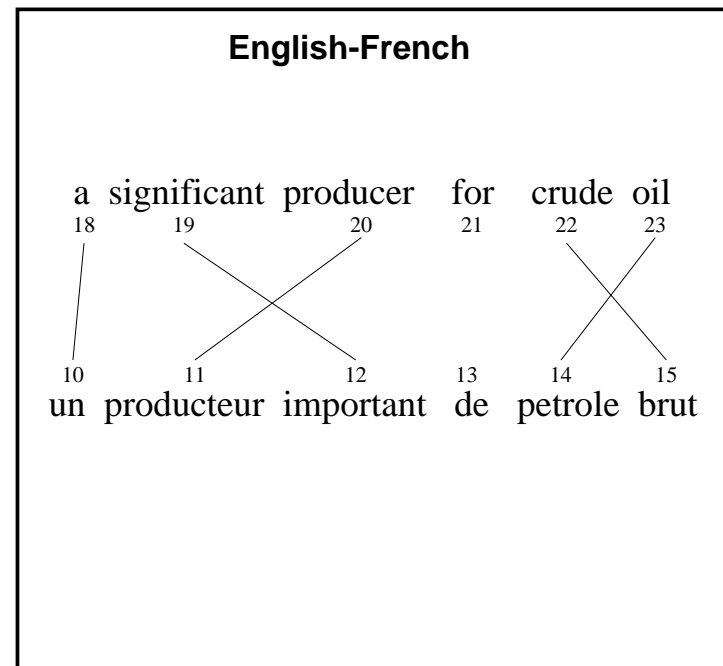
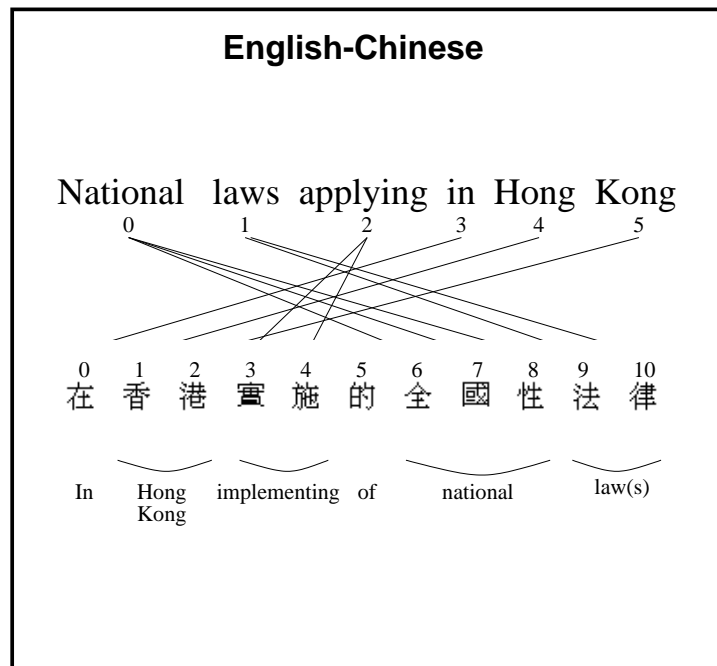
- Frequency and Context similarity models assume that words must start with the same initial letter at first iteration. This is not an unreasonable assumption for suffixal languages.

## Performance on Irregular Verb Sample

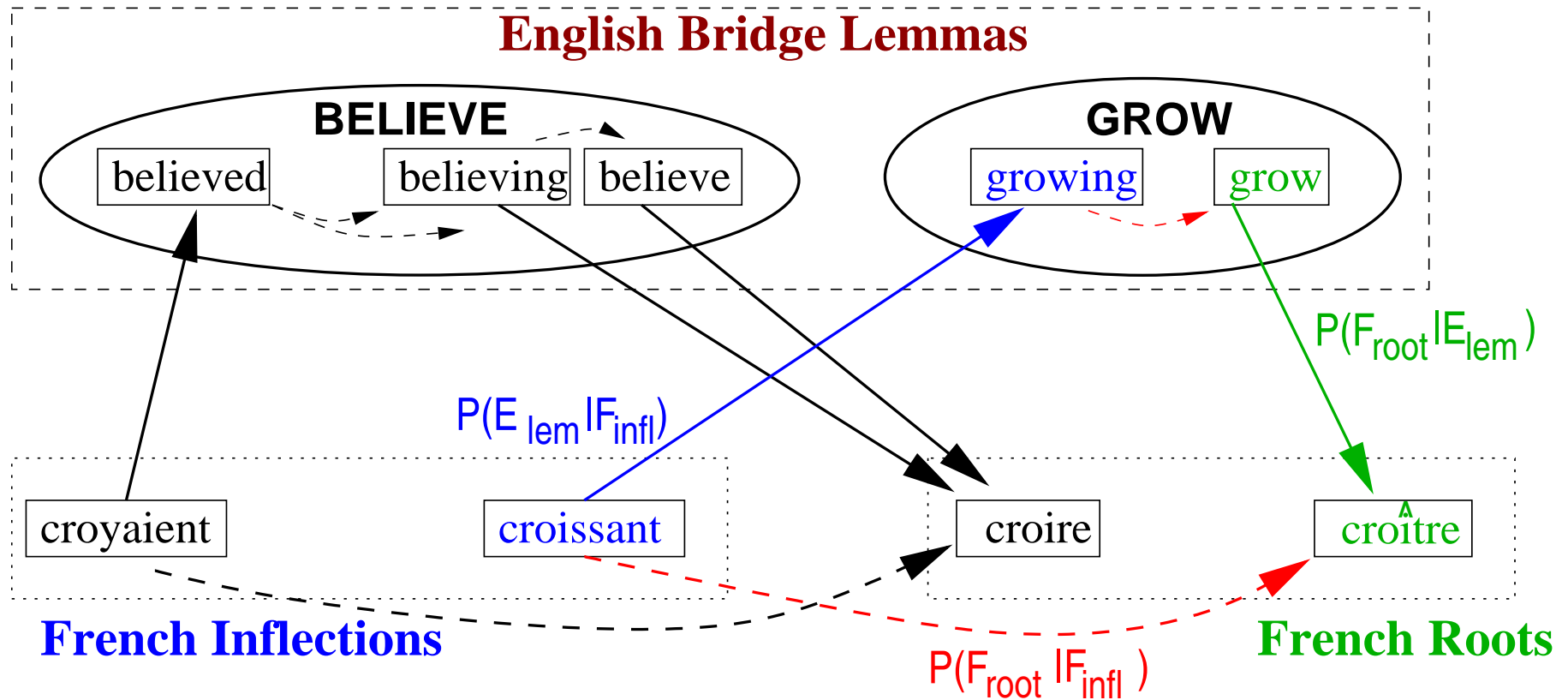
Word	True Root	CS+FS+LS+MS		(Itr 1)	CS+FS+LS (Itr 1)	CS+FS (Itr 1)	LS only (Itr 1)
		(Conv)g	Score				
got	get	go	1.30	go	go	go	gut
took	take	<b>take</b>	1.50	<b>take</b>	<b>take</b>	<b>take</b>	toot
became	become	<b>become</b>	2.35	<b>become</b>	<b>become</b>	<b>become</b>	<b>become</b>
clung	cling	<b>cling</b>	2.55	<b>cling</b>	<b>cling</b>	<b>cling</b>	<b>cling</b>
swore	swear	<b>swear</b>	2.80	<b>swear</b>	<b>swear</b>	<b>swear</b>	store
came	come	<b>come</b>	3.55	<b>come</b>	<b>come</b>	<b>come</b>	<b>come</b>
flung	fling	<b>fling</b>	4.60	<b>fling</b>	<b>fling</b>	<b>fling</b>	<b>fling</b>
strove	strive	<b>strive</b>	5.85	<b>strive</b>	<b>strive</b>	straddle	<b>strive</b>
swept	sweep	<b>sweep</b>	6.20	<b>sweep</b>	<b>sweep</b>	<b>sweep</b>	swap
woke	wake	<b>wake</b>	6.95	<b>wake</b>	<b>wake</b>	wind	<b>wake</b>
bore	bear	<b>bear</b>	7.75	<b>bear</b>	bar	<b>bear</b>	bare
lent	lend	<b>lend</b>	9.25	<b>lend</b>	<b>lend</b>	<b>lend</b>	<b>lend</b>
struck	strike	<b>strike</b>	11.60	<b>strike</b>	<b>strike</b>	<b>strike</b>	strut
bit	bite	<b>bite</b>	13.60	<b>bite</b>	<b>bite</b>	betray	bet
dove	dive	<b>dive</b>	17.25	<b>dive</b>	<b>dive</b>	dash	<b>dive</b>
caught	catch	<b>catch</b>	18.35	<b>catch</b>	cut	<b>catch</b>	cough
dealt	deal	<b>deal</b>	21.45	<b>deal</b>	<b>deal</b>	disagree	<b>deal</b>

## Multilingual Projection

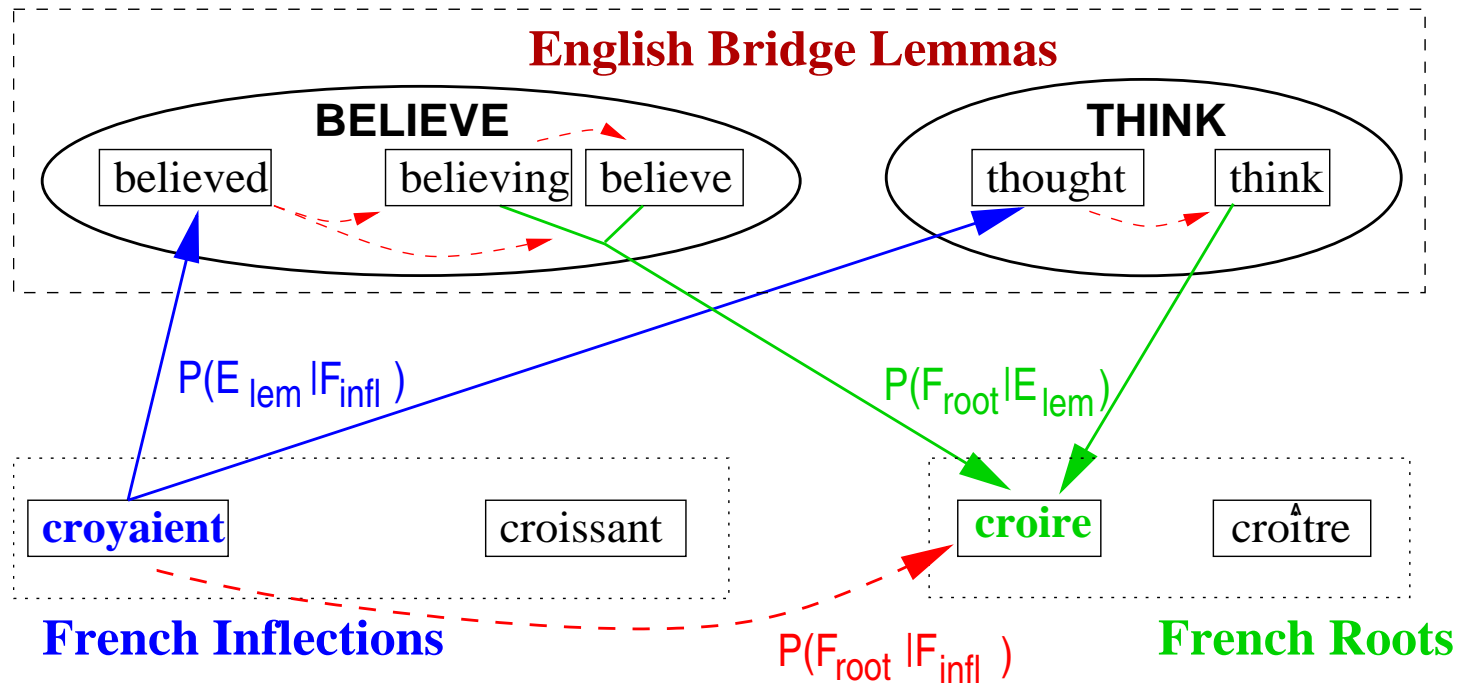
- Historically, large NLP investment in tools for English and a handful of other languages (French, Japanese)
- Use existing morphological analyzers to bootstrap training data for morphological analyzer for a second language
- Word align bilingual corpus (EGYPT, Y. Al-Onaizan et al. 1999):



# Morphological Analysis via Translingual Bridges: Leverage investment in existing analyzers



# General Alignment Model via Multiple Bridge Lemmas



$$P_{mp}(F_{root} | F_{infl}) = \sum_i P_a(F_{root} | E_{lem_i}) P_a(E_{lem_i} | F_{infl})$$

$$P_{mp}(\text{croire} | \text{croyaient}) = P(\text{croire} | \text{BELIEVE}) P(\text{BELIEVE} | \text{croyaient}) + P(\text{croire} | \text{THINK}) P(\text{THINK} | \text{croyaient}) + \dots$$

# Examples of Induced Morphological Analyses

## Induced Morphological Analyses for CZECH

Inflection	Root	POST Analysis	TopBridge
bral	brát	át→a +l	marry
brala	brát	át→a +la	accept
brali	brát	át→a +li	marry
byl	být	ýt→y +l	be
byli	být	ýt→y +li	be
bylo	být	ýt→y +lo	be
chovala	chovat	t→ε +la	behave
chová	chovat	at→ε +á	behave
chováme	chovat	at→ε +áme	behave
chodila	chodit	t→ε +la	walk
chodí	chodit	it→ε +í	walk
choďte	chodit	dit→ďt +e	swim
chránila	chránit	t→ε +la	protect
chrání	chránit	it→ε +í	protect
couval	couvat	t→ε +l	back
chcete	chtít	tít→c +ete	want
chceš	chtít	tít→c +eš	want
chci	chtít	tít→c +i	want
chtějí	chtít	ít→ěj +í	want
chtěli	chtít	ít→ě +li	want
chtělo	chtít	ít→ě +lo	want

## Induced Morphological Analyses for FRENCH

Inflection	Root	POST Analysis	TopBridge
abrège	abréger	éger→èg +e	shorten
abrègent	abréger	éger→èg +ent	shorten
abrègerai	abréger	er→ε +erai	curtail
achète	acheter	eter→èt +e	buy
achètent	acheter	eter→èt +ent	buy
achètera	acheter	eter→èt +era	buy
advieendrait	advenir	enir→iendr +ait	happen
advient	advenir	enir→ien +t	happen
aliène	aliéner	éner→èn +e	alienate
aliènent	aliéner	éner→èn +ent	alienate
conçu	concevoir	cevoir→ç +u	conceive
crois	croire	re→ε +s	believe
croyaient	croire	ire→y +aient	believe

## Induced Morphological Analyses for SPANISH

Inflection	Root	POST Analysis	TopBridge
aborreció	aborrecer	er→ε +ió	hate
aborrecía	aborrecer	er→ε +ía	hate
aborrezco	aborrecer	cer→zc +o	hate
abrace	abrazar	zar→c +e	embrace
abrazado	abrazar	ar→ε +ado	embrace
adquiere	adquirir	rir→er +e	get
anden	andar	ar→ε +en	walk
anduvo	andar	ar→uv +o	walk
buscáis	buscar	ar→ε +áis	seek
busque	buscar	car→qu +e	seek
busqué	buscar	car→qu +é	seek



## Use projected pairs as training for supervised models

Model	Precision		Coverage	
	Typ	Tok	Typ	Tok

### FRENCH Verbal Morphology Induction

French Hansards (12M words):

MProj only	.992	.999	.779	.994
MProj+POST	.998	.999	.988	.999
MProj+POST+BKM	<b>.994</b>	<b>.999</b>	1.00	1.00

French Hansards (1.2M words):

MProj only	.985	.998	.327	.976
MProj+POST	.995	.999	.958	.998
MProj+POST+BKM	<b>.979</b>	<b>.998</b>	1.00	1.00

French Hansards (120K words):

MProj only	.962	.931	.095	.901
MProj+POST	.984	.993	.916	.994
MProj+POST+BKM	<b>.932</b>	<b>.989</b>	1.00	1.00

French Bible (300K words) via English Bible:

MProj only	1.00	1.00	.052	.747
MProj+POST	.991	.998	.918	.992
MProj+POST+BKM	<b>.954</b>	<b>.994</b>	1.00	1.00

---

# Use projected pairs as training for supervised models (continued)

## CZECH Verbal Morphology Induction

Czech Reader's Digest (500K words):

MProj only	.915	.993	.152	.805
MProj+POST	.916	.917	.893	.975
MProj+POST+BKM	<b>.878</b>	<b>.913</b>	1.00	1.00

## SPANISH Verbal Morphology Induction

Spanish Bible (300K words) via English Bible:

MProj only	.973	.935	.264	.351
MProj+POST	.988	.998	.971	.967
MProj+POST+BKM	<b>.966</b>	<b>.985</b>	1.00	1.00

Spanish Bible (300K words) via French Bible:

MProj only	.980	.935	.722	.765
MProj+POST	.983	.974	.986	.993
MProj+POST+BKM	<b>.974</b>	<b>.968</b>	1.00	1.00

## Lemmatization Induction via Multiple Bible Versions

- Bible is easily alignable, available in many languages

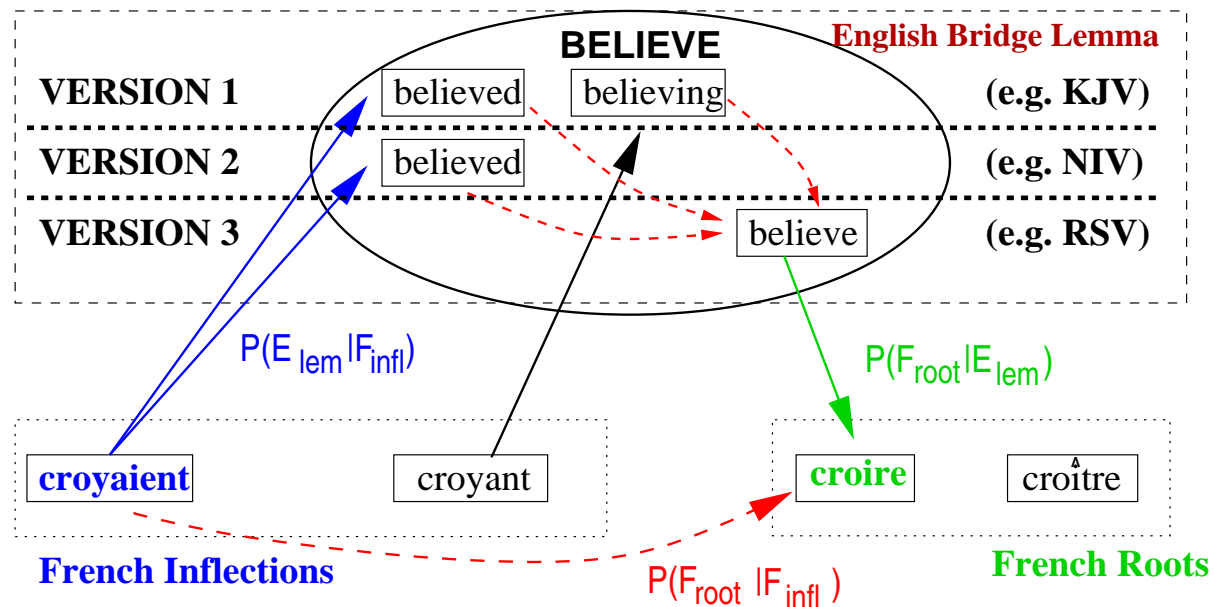
.954 accuracy (by type)
-------------------------

single French Bible

.994 accuracy (by token)
--------------------------

on full *modern French* test set

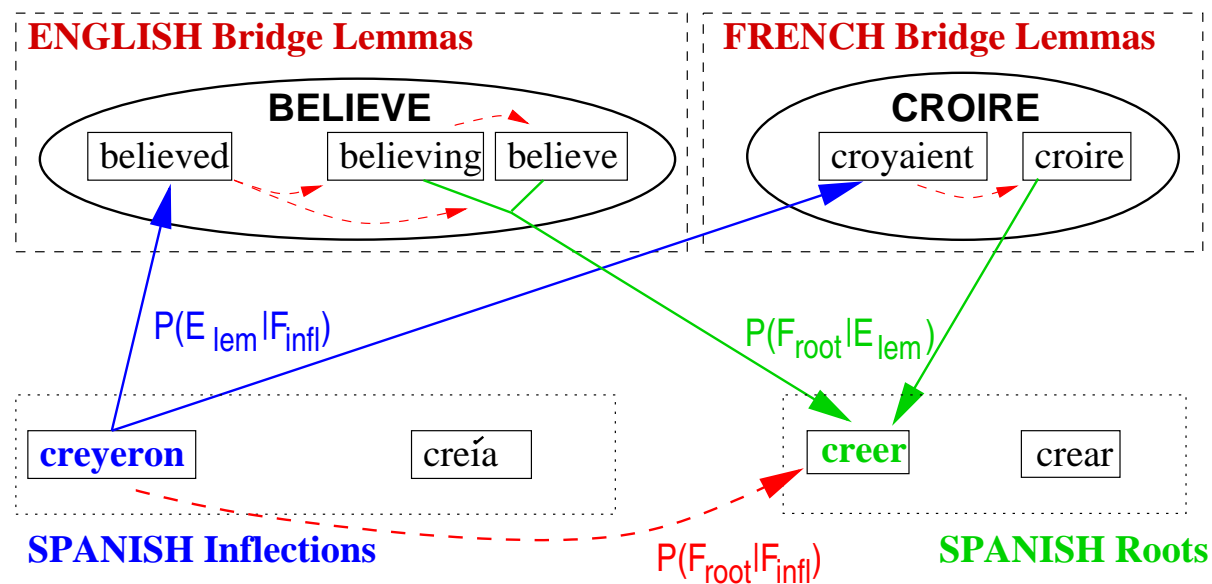
- Augment bridge potential using multiple English versions



⇒ .954 → .964 with *no* additional French resources

## Lemmatization Induction via Multiple Languages

- Using bitext in multiple languages adds bridging pathways
- Use newly lemmatized Bible in French to improve Spanish analysis



Spanish via 1 English Bible	.966	
Spanish via 1 French Bible	.974	(accuracy
Spanish via English + French Bibles	.981	by type)

---

## Use projected pairs as training for supervised models (continued)

Model	Precision		Coverage	
	Typ	Tok	Typ	Tok

### FRENCH Verbal Morphology Induction

French Bible (300K words) via 3 English Bibles:

MProj only	.928	.975	.100	.820
MProj+POST	.981	.991	.931	.990
MProj+POST+BKM	<b>.964</b>	<b>.991</b>	1.00	1.00

### SPANISH Verbal Morphology Induction

Spanish Bible (300K words) via 3 English Bibles:

MProj only	.964	.948	.468	.551
MProj+POST	.990	.998	.978	.987
MProj+POST	<b>.976</b>	<b>.987</b>	1.00	1.00

# Performance of Lemmatization Induction by Corpus Size

