

Statistical NLP

Spring 2008



Lecture 11: Word Alignment

Dan Klein – UC Berkeley

Machine Translation: Examples

Atlanta, preso il killer del palazzo di Giustizia

ATLANTA - La grande paura che per 26 ore ha attanagliato Atlanta è finita: Brian Nichols, l'uomo che aveva ucciso tre persone a palazzo di Giustizia e che ha poi ucciso un agente di dogana, s'è consegnato alla polizia, dopo avere cercato rifugio nell'alloggio di una donna in un complesso d'appartamenti alla periferia della città. Per tutto il giorno, il centro della città, sede della Coca Cola e dei Giochi 1996, cuore di una popolosa area metropolitana, era rimasto paralizzato.

Atlanta, taken the killer of the palace of Justice

ATLANTA - The great fear that for 26 hours has gripped Atlanta is ended: Brian Nichols, the man who had killed three persons to palace of Justice and that a customs agent has then killed, s' is delivered to the police, after to have tried shelter in the lodging of one woman in a complex of apartments to the periphery of the city. For all the day, the center of the city, center of the Coke Strains land of Giochi 1996, heart of one popolosa metropolitan area, was remained paralyzed.

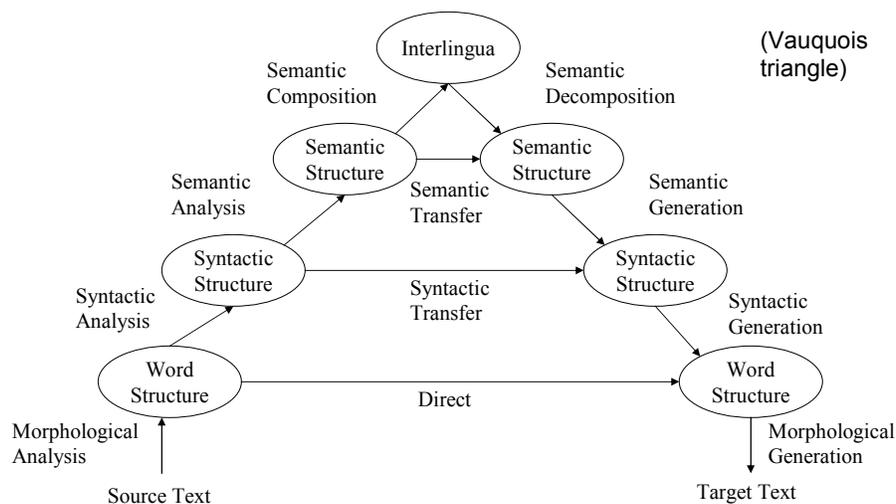
Machine Translation

Madame la présidente, votre présidence de cette institution a été marquante.
 Mrs Fontaine, your presidency of this institution has been outstanding.
 Madam President, president of this house has been discoveries.
 Madam President, your presidency of this institution has been impressive.

Je vais maintenant m'exprimer brièvement en irlandais.
 I shall now speak briefly in Irish .
 I will now speak briefly in Ireland .
 I will now speak briefly in Irish .

Nous trouvons en vous un président tel que nous le souhaitons.
 We think that you are the type of president that we want.
 We are in you a president as the wanted.
 We are in you a president as we the wanted.

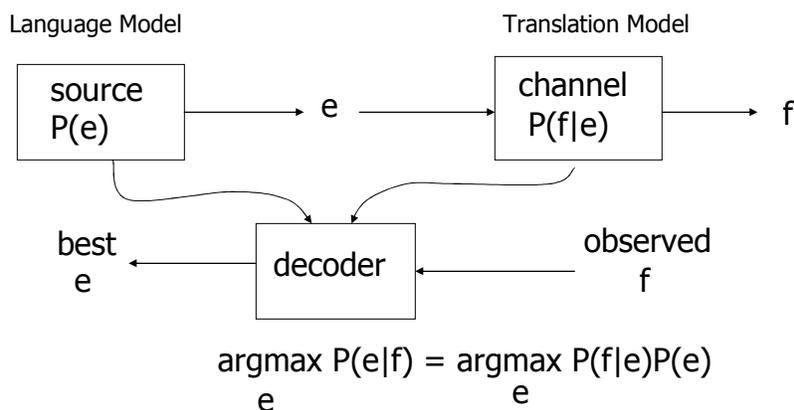
Levels of Transfer



General Approaches

- **Rule-based approaches**
 - Expert system-like rewrite systems
 - Interlingua methods (analyze and generate)
 - Lexicons come from humans
 - Can be very fast, and can accumulate a lot of knowledge over time (e.g. Systran)
- **Statistical approaches**
 - Word-to-word translation
 - Phrase-based translation
 - Syntax-based translation (tree-to-tree, tree-to-string)
 - Trained on parallel corpora
 - Usually noisy-channel (at least in spirit)

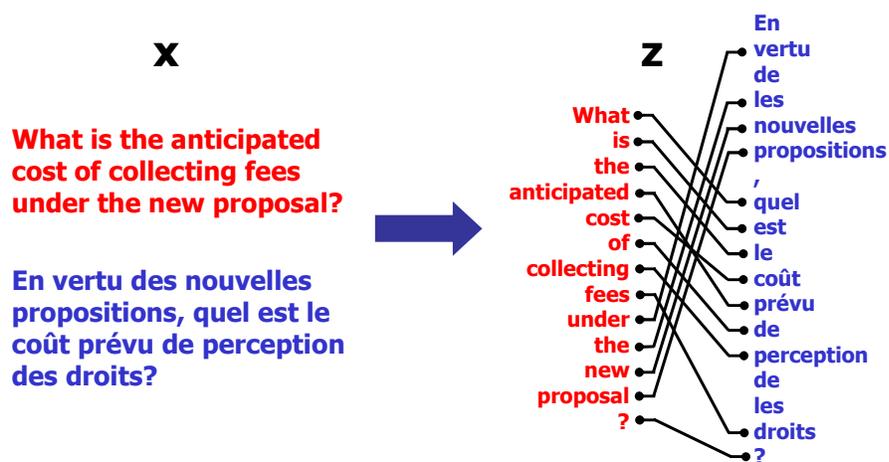
MT System Components



Today

- The components of a simple MT system
 - You already know about the LM
 - Word-alignment based TMs
 - IBM models 1 and 2, HMM model
 - A simple decoder
- Next few classes
 - More complex word-level and phrase-level TMs
 - Tree-to-tree and tree-to-string TMs
 - More sophisticated decoders

Word Alignment



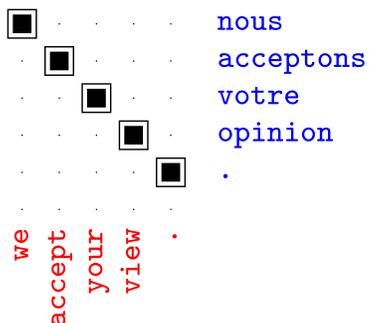
Unsupervised Word Alignment

- Input: a *bitext*: pairs of translated sentences

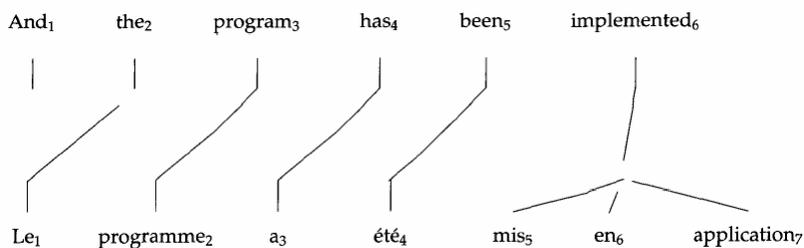
nous	acceptons	votre	opinion	.
we	accept	your	view	.

- Output: *alignments*: pairs of translated words

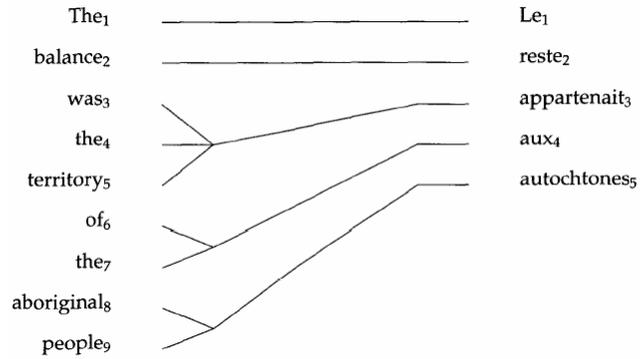
- When words have unique sources, can represent as a (forward) alignment function a from French to English positions



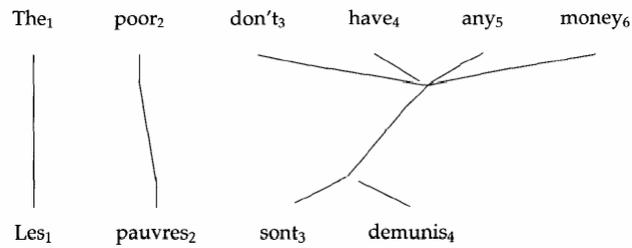
1-to-Many Alignments



Many-to-1 Alignments



Many-to-Many Alignments



A Word-Level TM?

- What might a model of $P(f|e)$ look like?

$e = e_1 \dots e_I$ And₁ the₂ program₃ has₄ been₅ implemented₆
 $f = f_1 \dots f_J$ Le₁ programme₂ a₃ été₄ mis₅ en₆ application₇

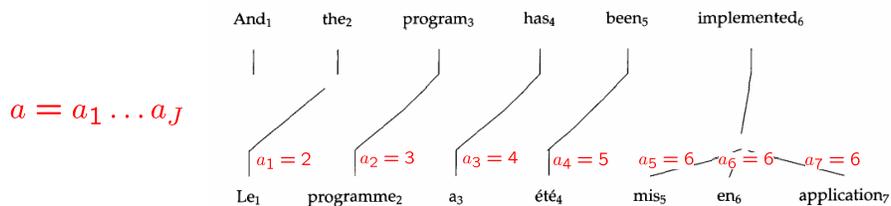
$$P(f|e) = \prod_j P(f_j | e_1 \dots e_I)$$

How to estimate this?

What can go wrong here?

IBM Model 1 (Brown 93)

- Alignments: a hidden vector called an *alignment* specifies which English source is responsible for each French target word.



$$\begin{aligned}
 P(f, a|e) &= \prod_j P(a_j = i) P(f_j | e_i) \\
 &= \prod_j \frac{1}{I+1} P(f_j | e_i)
 \end{aligned}$$

$$P(f|e) = \sum_a P(f, a|e)$$

IBM Model 1

- Obvious first stab: greedy matchings
- Better approach: re-estimated generative models

$$P(f|e) = \sum_a P(f, a|e)$$

$$P(f, a|e) = \prod_j P(a_j = i|e) P(f_j|e_i)$$

$$P(a_j = i|e, f) = \frac{P(f_j|e_i)}{\sum_{i'} P(f_j|e_{i'})}$$

- Basic idea: pick a source for each word, update co-occurrence statistics, repeat

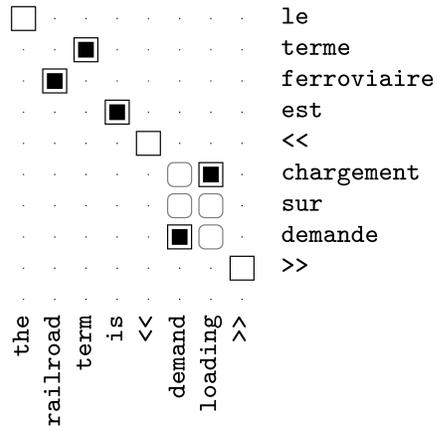
Evaluating TMs

- How do we measure quality of a word-to-word model?
 - Method 1: use in an end-to-end translation system
 - Hard to measure translation quality
 - Option: human judges
 - Option: reference translations (NIST, BLEU)
 - Option: combinations (HTER)
 - Actually, no one uses word-to-word models alone as TMs
 - Method 2: measure quality of the alignments produced
 - Easy to measure
 - Hard to know what the gold alignments should be
 - Often does not correlate well with translation quality (like perplexity in LMs)

Intersected Model 1

- Post-intersection: standard practice to train models in each direction then intersect their predictions [Och and Ney, 03]
- Second model is basically a filter on the first
 - Precision jumps, recall drops
 - End up not guessing hard alignments

Model	P/R	AER
Model 1 E→F	82/58	30.6
Model 1 F→E	85/58	28.7
Model 1 AND	96/46	34.8



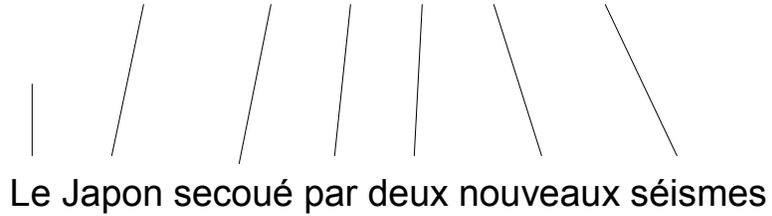
Joint Training?

- Overall:
 - Similar high precision to post-intersection
 - But recall is much higher
 - More confident about positing non-null alignments

Model	P/R	AER
Model 1 E→F	82/58	30.6
Model 1 F→E	85/58	28.7
Model 1 AND	96/46	34.8
Model 1 INT	93/69	19.5

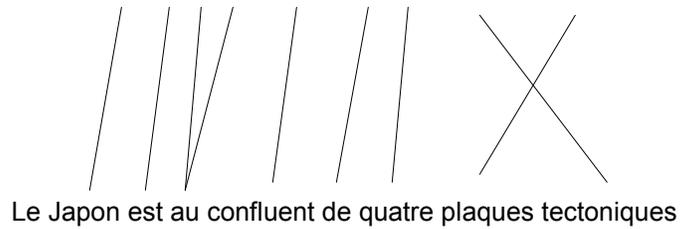
Monotonic Translation

Japan shaken by two new quakes



Local Order Change

Japan is at the junction of four tectonic plates



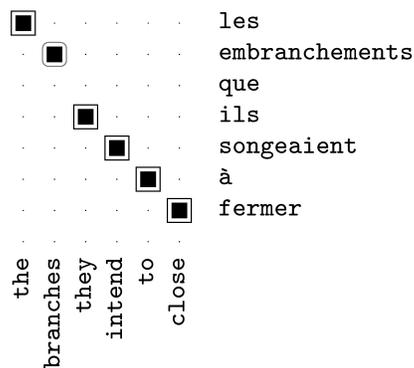
IBM Model 2

- Alignments tend to the diagonal (broadly at least)

$$P(f, a|e) = \prod_j P(a_j = i|j, I, J) P(f_j|e_i)$$
$$P(\text{dist} = i - j \frac{I}{J})$$
$$\frac{1}{Z} e^{-\alpha(i - j \frac{I}{J})}$$

- Other schemes for biasing alignments towards the diagonal:
 - Relative vs absolute alignment
 - Asymmetric distances
 - Learning a full multinomial over distances

Example



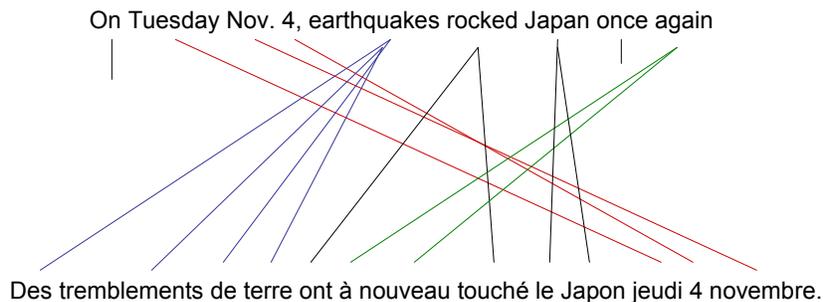
EM for Models 1/2

- Model 1 Parameters:
 - Translation probabilities (1+2) $P(f_j|e_i)$
 - Distortion parameters (2 only) $P(a_j = i|j, I, J)$
- Start with $P(f_j|e_i)$ uniform, including $P(f_j|null)$
- For each sentence:
 - For each French position j
 - Calculate posterior over English positions

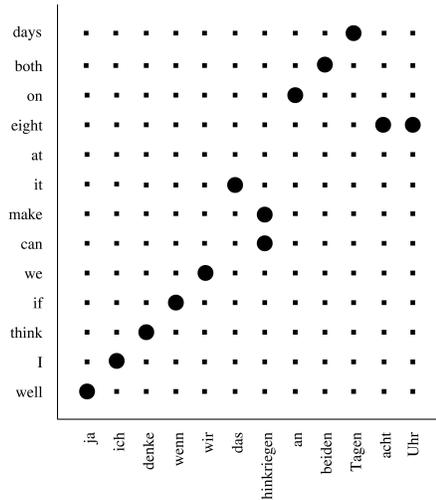
$$P(a_j = i|f, e) = \frac{P(a_j = i|j, I, J)P(f_j|e_i)}{\sum_{i'} P(a_j = i'|j, I, J)P(f_j|e_{i'})}$$

- (or just use best single alignment)
 - Increment count of word f_j with word e_i by these amounts
 - Also re-estimate distortion probabilities for model 2
- Iterate until convergence

Phrase Movement



Phrase Movement



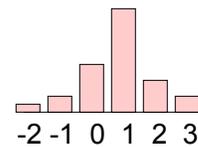
The HMM Model

- Model 2 preferred global monotonicity
- We want local monotonicity:
 - Most jumps are small
- HMM model (Vogel 96)

f	t(f e)
nationale	0.469
national	0.418
nationaux	0.054
nationales	0.029

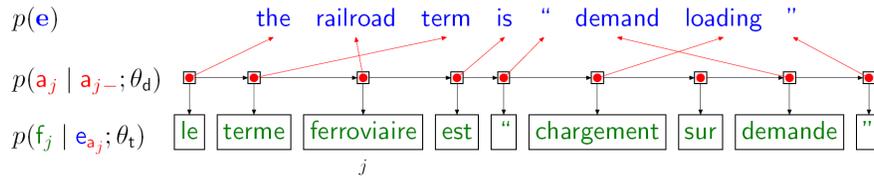
$$P(f, a|e) = \prod_j P(a_j|a_{j-1})P(f_j|e_i)$$

$$P(a_j - a_{j-1})$$



- Re-estimate using the forward-backward algorithm
- Handling nulls requires some care
- What are we still missing?

The HMM Model



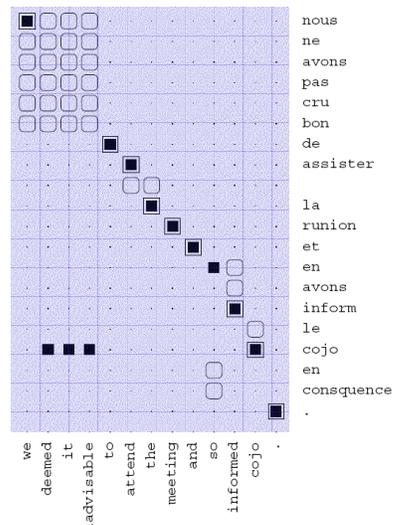
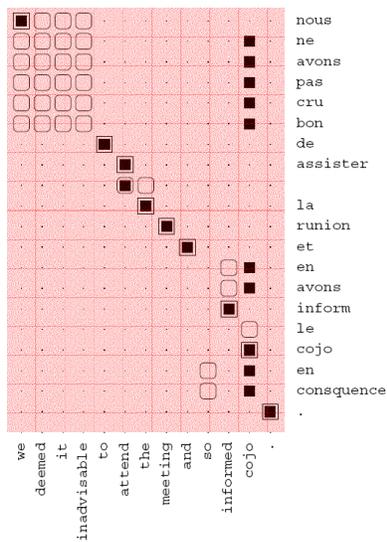
Distortion θ_d

$$\begin{aligned}
 p(\uparrow \uparrow) &= 0.6 \\
 p(\uparrow \rightarrow) &= 0.2 \\
 p(\rightarrow \times) &= 0.1 \\
 &\dots
 \end{aligned}$$

Translation θ_t

$$\begin{aligned}
 p(\text{the} \rightarrow \text{le}) &= 0.53 \\
 p(\text{the} \rightarrow \text{la}) &= 0.24 \\
 p(\text{railroad} \rightarrow \text{ferroviaire}) &= \mathbf{0.19} \\
 p(\text{NULL} \rightarrow \text{le}) &= 0.12 \\
 &\dots
 \end{aligned}$$

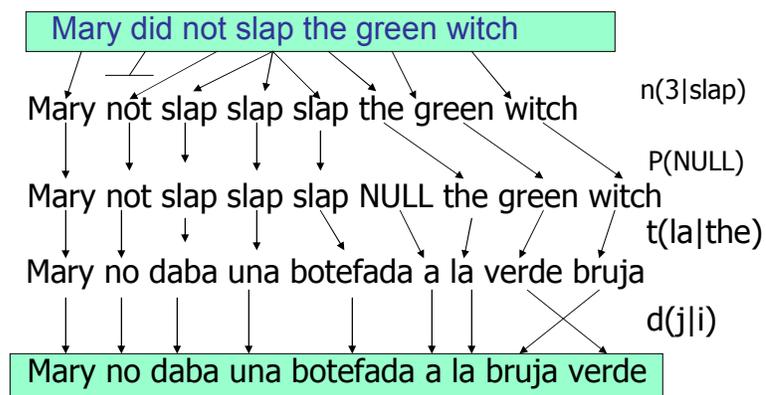
HMM Examples



AER for HMMs

Model	AER
Model 1 INT	19.5
HMM E→F	11.4
HMM F→E	10.8
HMM AND	7.1
HMM INT	4.7
GIZA M4 AND	6.9

IBM Models 3/4/5



[from Al-Onaizan and Knight, 1998]

Examples: Translation and Fertility

the

f	$t(f e)$	ϕ	$n(\phi e)$
le	0.497	1	0.746
la	0.207	0	0.254
les	0.155		
l'	0.086		
ce	0.018		
cette	0.011		

not

f	$t(f e)$	ϕ	$n(\phi e)$
ne	0.497	2	0.735
pas	0.442	0	0.154
non	0.029	1	0.107
rien	0.011		

farmers

f	$t(f e)$	ϕ	$n(\phi e)$
agriculteurs	0.442	2	0.731
les	0.418	1	0.228
cultivateurs	0.046	0	0.039
producteurs	0.021		

Example: Idioms

nodding

f	$t(f e)$	ϕ	$n(\phi e)$
signe	0.164	4	0.342
la	0.123	3	0.293
tête	0.097	2	0.167
oui	0.086	1	0.163
fait	0.073	0	0.023
que	0.073		
hoche	0.054		
hocher	0.048		
faire	0.030		
me	0.024		
approuve	0.019		
qui	0.019		
un	0.012		
faites	0.011		

Example: Morphology

should

f	$t(f e)$	ϕ	$n(\phi e)$
devrait	0.330	1	0.649
devraient	0.123	0	0.336
devrions	0.109	2	0.014
faudrait	0.073		
faut	0.058		
doit	0.058		
aurait	0.041		
doivent	0.024		
devons	0.017		
devrais	0.013		

Some Results

- [Och and Ney 03]

Model	Training scheme	0.5K	8K	128K	1.47M
Dice		50.9	43.4	39.6	38.9
Dice+C		46.3	37.6	35.0	34.0
Model 1	1^5	40.6	33.6	28.6	25.9
Model 2	$1^5 2^5$	46.7	29.3	22.0	19.5
HMM	$1^5 H^5$	26.3	23.3	15.0	10.8
Model 3	$1^5 2^5 3^3$	43.6	27.5	20.5	18.0
	$1^5 H^5 3^3$	27.5	22.5	16.6	13.2
Model 4	$1^5 2^5 3^3 4^3$	41.7	25.1	17.3	14.1
	$1^5 H^5 3^3 4^3$	26.1	20.2	13.1	9.4
	$1^5 H^5 4^3$	26.3	21.8	13.3	9.3
Model 5	$1^5 H^5 4^3 5^3$	26.5	21.5	13.7	9.6
	$1^5 H^5 3^3 4^3 5^3$	26.5	20.4	13.4	9.4
Model 6	$1^5 H^5 4^3 6^3$	26.0	21.6	12.8	8.8
	$1^5 H^5 3^3 4^3 6^3$	25.9	20.3	12.5	8.7

Stack Decoding

- **Stack decoding:**
 - Beam search
 - Usually A* estimates for completion cost
 - One stack per candidate sentence length
- **Other methods:**
 - Dynamic programming decoders possible if we make assumptions about the set of allowable permutations

sent length	decoder type	time (sec/sent)	search errors	translation errors (semantic and/or syntactic)	NE	PME	DSE	FSE	HSE	CE
6	IP	47.50	0	57	44	57	0	0	0	0
6	stack	0.79	5	58	43	53	1	0	0	4
6	greedy	0.07	18	60	38	45	5	2	1	10
8	IP	499.00	0	76	27	74	0	0	0	0
8	stack	5.67	20	75	24	57	1	2	2	15
8	greedy	2.66	43	75	20	38	4	5	1	33