

Introduction to Statistical Machine Translation

David Kauchak
CS159 – Fall 2014

Some slides adapted from

Philipp Koehn
School of Informatics
University of Edinburgh

Kevin Knight
USC/Information Sciences Institute
USC/Computer Science Department

Dan Klein
Computer Science Department
UC Berkeley

Admin

Assignment 5

Quiz #2

- Mean 26.5 (88%)
- Median 26.5

Language translation



MT Systems

Where have you seen machine translation systems?

Machine Translation

美国关岛国际机场及其办公室均接获一名自称沙地阿拉伯富商拉登等发出的电子邮件，威胁将会向机场等公众地方发动生化袭击後，关岛经保持高度戒备。



The U.S. island of Guam is maintaining a high state of alert after the Guam airport and its offices both received an e-mail from someone calling himself the Saudi Arabian Osama bin Laden and threatening a biological/chemical attack against public places such as the airport .

The classic acid test for natural language processing.

Requires capabilities in both interpretation and generation.

“People around the world stubbornly refuse to write everything in English.” ☺

Machine Translation

美国关岛国际机场及其办公室均接获一名自称沙地阿拉伯富商拉登等发出的电子邮件，威胁将会向机场等公众地方发动生化袭击後，关岛经保持高度戒备。

Machine translation is becoming very prevalent

Even PowerPoint has translation built into it!

The U.S. island of Guam is maintaining a high state of alert after the Guam airport and its offices both received an e-mail from someone calling himself the Saudi Arabian Osama bin Laden and threatening a biological/chemical attack against public places such as the airport .

United States Office of the Guam International Airport and were received by a man claiming to be Saudi Arabian businessman Osama bin Laden, sent emails, threats to airports and other public places will launch biological or chemical attack, remain on high alert in Guam.

Warren Weaver (1947)



ingcmpnqsnwf cv fpn owoktvcv

hu ihgzsnwfv rqcffnw cw owgcnwf

kowazoanv ...

Warren Weaver (1947)



e e e e
ingcmpnqsnwf cv fpn owoktvcv

e e e
hu ihgzsnwfv rqcffnw cw owgcnwf

e
kowazoanv ...

Warren Weaver (1947)



e e e the
ingcmpnqsnwf cv fpn owoktvcv
e e e
hu ihgzsnwfv rqcffnw cw owgcnwf
e
kowazoanv ...

Warren Weaver (1947)



e he e the
ingcmpnqsnwf cv fpn owoktvcv
e e e t
hu ihgzsnwfv rqcffnw cw owgcnwf
e
kowazoanv ...

Warren Weaver (1947)



e he e of the
ingcmpnqsnwf cv fpn owoktvcv
e e e t
hu ihgzsnwfv rqcffnw cw owgcnwf
e
kowazoanv ...

Warren Weaver (1947)



e he e of the fof
ingcmpnqsnwf cv fpn owoktvcv
e f o e o o e t
hu ihgzsnwfv rqcffnw cw owgcnwf
ef
kowazoanv ...

Warren Weaver (1947)



e he e ~~o~~ the
ingcmpnqsnwf cv fpn owoktvcv
e e e t
hu ihgzsnwfv rqcffnw cw owgcnwf
e
kowazoanv ...

Warren Weaver (1947)



e he e is the sis
ingcmpnqsnwf cv fpn owoktvcv
e s i e i i e t
hu ihgzsnwfv rqcffnw cw owgcnwf
es
kowazoanv ...

Warren Weaver (1947)



decipherment is the analysis
ingcmpnqsnwf cv fpn owoktvcv
of documents written in ancient
hu ihgzsnwfv rqcffnw cw owgcnwf
languages ...
kowazoanv ...

Warren Weaver (1947)

Can this be computerized?

The non-Turkish guy next to me is even deciphering Turkish! All he needs is a statistical table of letter-pair frequencies in Turkish ...



Collected mechanically from a Turkish body of text, or *corpus*



“When I look at an article in Russian, I say: this is really written in English, but it has been coded in some strange symbols. I will now proceed to decode.”
 - Warren Weaver, March 1947

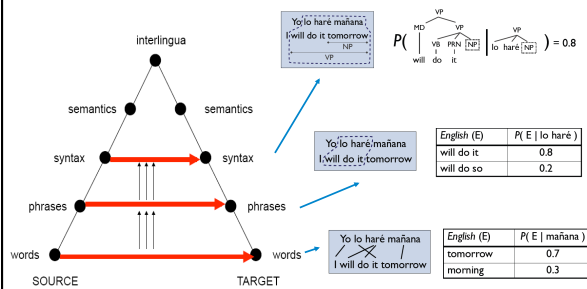


“When I look at an article in Russian, I say: this is really written in English, but it has been coded in some strange symbols. I will now proceed to decode.”
 - Warren Weaver, March 1947



“... as to the problem of mechanical translation, I frankly am afraid that the [semantic] boundaries of words in different languages are too vague ... to make any quasi-mechanical translation scheme very hopeful.”
 - Norbert Wiener, April 1947

Levels of Transfer



World-Level MT: Examples

la politique de la haine . (Foreign Original)
 politics of hate . (Reference Translation)
 the policy of the hatred . (IBM4+N-grams+Stack)

nous avons signé le protocole . (Foreign Original)
 we did sign the memorandum of agreement . (Reference Translation)
 we have signed the protocol . (IBM4+N-grams+Stack)

où était le plan solide ? (Foreign Original)
 but where was the solid plan ? (Reference Translation)
 where was the economic base ? (IBM4+N-grams+Stack)

Phrasal / Syntactic MT: Examples

Le président américain Barack Obama doit annoncer lundi de nouvelles mesures en faveur des constructeurs automobile. General motors et Chrysler avaient déjà bénéficié fin 2008 d'un prêt d'urgence cumulé de 17,4 milliards de dollars, et ont soumis en février au Trésor un plan de restructuration basé sur un total de 22 milliards de dollars d'aides publiques supplémentaires.

Interrogé sur la chaîne CBS dimanche, le président a toutefois clairement précisé que le gouvernement ne prêterait pas d'argent sans de fortes contreparties. "Il faudra faire des sacrifices à tous les niveaux", a-t-il prévenu. "Tout le monde devra se réunir autour de la table et se mettre d'accord sur une restructuration en profondeur".

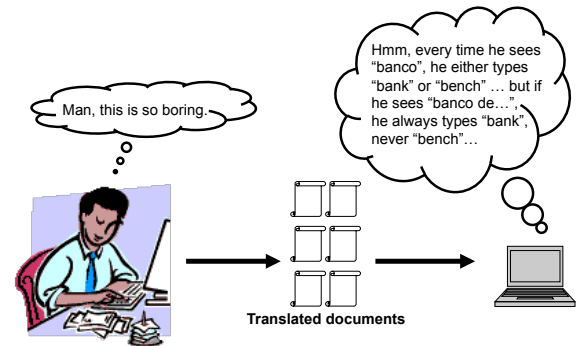
General Motors et Chrysler sont engagés dans des négociations avec le principal syndicat de l'automobile. Les constructeurs souhaitent diminuer leurs cotisations aux caisses de retraites, et accorder en échange des actions aux syndicats. Ils souhaiteraient également négocier des baisses des salaires.

U.S. President Barack Obama to announce Monday new measures to help automakers. General Motors and Chrysler had already received late in 2008 a cumulative emergency loan of 17.4 billion dollars, and submitted to the Treasury in February in a restructuring plan based on a total of 22 billion dollars in additional aid.

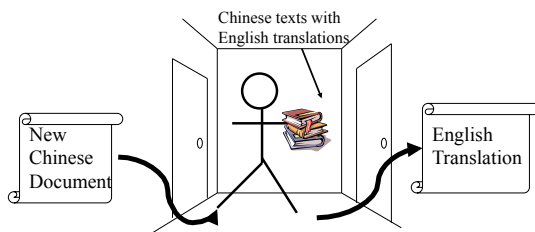
Interviewed on CBS Sunday, the president has clearly stated that the government does not lend money without strong counterparts. "We must make sacrifices at all levels," he warned. "Everyone should gather around the table and agree on a profound restructuring."

General Motors and Chrysler are engaged in negotiations with the major union of the car. Manufacturers wishing to reduce their contributions to pension funds, and give in exchange for the shares to trade unions. They would also negotiate lower wages.

Data-Driven Machine Translation



Welcome to the Chinese Room



You can teach yourself to translate Chinese using *only* bilingual data (without grammar books, dictionaries, any people to answer your questions...)

Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan: farok errok hihok yorok klok kantok ok-yurp

1a. ok-voon ororok sprok .	7a. lalok farok ororok lalok sprok izok enemok .
1b. at-voon bichat dat .	7b. wat jiat bichat wat dat vat enecat .
2a. ok-drubel ok-voon anak plok sprok .	8a. lalok brok anak plok nok .
2b. at-drubel at-voon pippat rrat dat .	8b. iat lat pippat rrat nnat .
3a. erok sprok izok hihok ghrok .	9a. wiwok nok izok kantok ok-yurp .
3b. totat dat arrat vat hilat .	9b. totat nnat quat oloat at-yurp .
4a. ok-voon anak drok brok jok .	10a. lalok mok nok yorok ghrok klok .
4b. at-voon krat pippat sat lat .	10b. wat nnat gat mat bat hilat .
5a. wiwok farok izok stok .	11a. lalok nok errok hihok yorok zanzanok .
5b. totat jiat quat cat .	11b. wat nnat arrat mat zanzanat .
6a. lalok sprok izok jok stok .	12a. lalok rarok nok izok hihok mok .
6b. wat dat krat quat cat .	12b. wat nnat forat arrat vat gat .

Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan: **farok** crrrok hihok yorok klok kantok ok-yurp

1a. ok-voon ororok sprok .	7a. lalok farok ororok lalok sprok izok enemok .
1b. at-voon bichat dat .	7b. wat jjat bichat wat dat vat eneat .
2a. ok-drubel ok-voon anak plok sprok .	8a. lalok brok anak plok nok .
2b. at-drubel at-voon pippat rrat dat .	8b. iat lat pippat rrat nnat .
3a. erok sprok izok hihok ghirok .	9a. wiwok nok izok kantok ok-yurp .
3b. totat dat arrat vat hilat .	9b. totat nnat quat oloat at-yurp .
4a. ok-voon anak drok brok jok .	10a. lalok mok nok yorok ghirok klok .
4b. at-voon krat pippat sat lat .	10b. wat nnat gat mat bat hilat .
5a. wiwok farok izok stok .	11a. lalok nok crrrok hihok yorok zanzanok .
5b. totat jjat quat cat .	11b. wat nnat arrat mat zanzanat .
6a. lalok sprok izok jok stok .	12a. lalok rarok nok izok hihok mok .
6b. wat dat krat quat cat .	12b. wat nnat forat arrat vat gat .

Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan: **farok** crrrok hihok yorok klok kantok ok-yurp

1a. ok-voon ororok sprok .	7a. lalok farok ororok lalok sprok izok enemok .
1b. at-voon bichat dat .	7b. wat jjat bichat wat dat vat eneat .
2a. ok-drubel ok-voon anak plok sprok .	8a. lalok brok anak plok nok .
2b. at-drubel at-voon pippat rrat dat .	8b. iat lat pippat rrat nnat .
3a. erok sprok izok hihok ghirok .	9a. wiwok nok izok kantok ok-yurp .
3b. totat dat arrat vat hilat .	9b. totat nnat quat oloat at-yurp .
4a. ok-voon anak drok brok jok .	10a. lalok mok nok yorok ghirok klok .
4b. at-voon krat pippat sat lat .	10b. wat nnat gat mat bat hilat .
5a. wiwok farok izok stok .	11a. lalok nok crrrok hihok yorok zanzanok .
5b. totat jjat quat cat .	11b. wat nnat arrat mat zanzanat .
6a. lalok sprok izok jok stok .	12a. lalok rarok nok izok hihok mok .
6b. wat dat krat quat cat .	12b. wat nnat forat arrat vat gat .

Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan: **farok** crrrok hihok yorok klok kantok ok-yurp

1a. ok-voon ororok sprok .	7a. lalok farok ororok lalok sprok izok enemok .
1b. at-voon bichat dat .	7b. wat jjat bichat wat dat vat eneat .
2a. ok-drubel ok-voon anak plok sprok .	8a. lalok brok anak plok nok .
2b. at-drubel at-voon pippat rrat dat .	8b. iat lat pippat rrat nnat .
3a. erok sprok izok hihok ghirok .	9a. wiwok nok izok kantok ok-yurp .
3b. totat dat arrat vat hilat .	9b. totat nnat quat oloat at-yurp .
4a. ok-voon anak drok brok jok .	10a. lalok mok nok yorok ghirok klok .
4b. at-voon krat pippat sat lat .	10b. wat nnat gat mat bat hilat .
5a. wiwok farok izok stok .	11a. lalok nok crrrok hihok yorok zanzanok .
5b. totat jjat quat cat .	11b. wat nnat arrat mat zanzanat .
6a. lalok sprok izok jok stok .	12a. lalok rarok nok izok hihok mok .
6b. wat dat krat quat cat .	12b. wat nnat forat arrat vat gat .

Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan: **farok** crrrok **hihok** yorok klok kantok ok-yurp

1a. ok-voon ororok sprok .	7a. lalok farok ororok lalok sprok izok enemok .
1b. at-voon bichat dat .	7b. wat jjat bichat wat dat vat eneat .
2a. ok-drubel ok-voon anak plok sprok .	8a. lalok brok anak plok nok .
2b. at-drubel at-voon pippat rrat dat .	8b. iat lat pippat rrat nnat .
3a. erok sprok izok hihok ghirok .	9a. wiwok nok izok kantok ok-yurp .
3b. totat dat arrat vat hilat .	9b. totat nnat quat oloat at-yurp .
4a. ok-voon anak drok brok jok .	10a. lalok mok nok yorok ghirok klok .
4b. at-voon krat pippat sat lat .	10b. wat nnat gat mat bat hilat .
5a. wiwok farok izok stok .	11a. lalok nok crrrok hihok yorok zanzanok .
5b. totat jjat quat cat .	11b. wat nnat arrat mat zanzanat .
6a. lalok sprok izok jok stok .	12a. lalok rarok nok izok hihok mok .
6b. wat dat krat quat cat .	12b. wat nnat forat arrat vat gat .

Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan: **farok** **errrok** **hihok** **yorok** **clok** kantok ok-yurp

1a. ok-voon ororok sprok .	7a. lalok farok ororok lalok sprok izok enemok .
1b. at-voon bichat dat .	7b. wat jjat bichat wat dat vat eneat .
2a. ok-drubel ok-voon anak plok sprok .	8a. lalok brok anak plok nok .
2b. at-drubel at-voon pippat rrat dat .	8b. iat lat pippat rrat nnat .
3a. erok sprok izok hihok ghirok .	9a. wiwok nok izok kantok ok-yurp .
3b. totat dat arrat vat hilat .	9b. totat nnat quat oloat at-yurp .
4a. ok-voon anak drok brok jok .	10a. lalok mok nok yorok ghirok clok .
4b. at-voon krat pippat sat lat .	10b. wat nnat gat mat bat hilat .
5a. wiwok farok izok stok .	11a. lalok nok errrok hihok yorok zanzanok .
5b. totat jjat quat cat .	11b. wat nnat arrat mat zanzanat .
6a. lalok sprok izok jok stok .	12a. lalok rarok nok izok hihok mok .
6b. wat dat krat quat cat .	12b. wat nnat forat arrat vat gat .

Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan: **farok** **errrok** **hihok** **yorok** **clok** kantok ok-yurp

1a. ok-voon ororok sprok .	7a. lalok farok ororok lalok sprok izok enemok .
1b. at-voon bichat dat .	7b. wat jjat bichat wat dat vat eneat .
2a. ok-drubel ok-voon anak plok sprok .	8a. lalok brok anak plok nok .
2b. at-drubel at-voon pippat rrat dat .	8b. iat lat pippat rrat nnat .
3a. erok sprok izok hihok ghirok .	9a. wiwok nok izok kantok ok-yurp .
3b. totat dat arrat vat hilat .	9b. totat nnat quat oloat at-yurp .
4a. ok-voon anak drok brok jok .	10a. lalok mok nok yorok ghirok clok .
4b. at-voon krat pippat sat lat .	10b. wat nnat gat mat bat hilat .
5a. wiwok farok izok stok .	11a. lalok nok errrok hihok yorok zanzanok .
5b. totat jjat quat cat .	11b. wat nnat arrat mat zanzanat .
6a. lalok sprok izok jok stok .	12a. lalok rarok nok izok hihok mok .
6b. wat dat krat quat cat .	12b. wat nnat forat arrat vat gat .

Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan: **farok** **errrok** **hihok** **yorok** **clok** kantok ok-yurp

1a. ok-voon ororok sprok .	7a. lalok farok ororok lalok sprok izok enemok .
1b. at-voon bichat dat .	7b. wat jjat bichat wat dat vat eneat .
2a. ok-drubel ok-voon anak plok sprok .	8a. lalok brok anak plok nok .
2b. at-drubel at-voon pippat rrat dat .	8b. iat lat pippat rrat nnat .
3a. erok sprok izok hihok ghirok .	9a. wiwok nok izok kantok ok-yurp .
3b. totat dat arrat vat hilat .	9b. totat nnat quat oloat at-yurp .
4a. ok-voon anak drok brok jok .	10a. lalok mok nok yorok ghirok clok .
4b. at-voon krat pippat sat lat .	10b. wat nnat gat mat bat hilat .
5a. wiwok farok izok stok .	11a. lalok nok errrok hihok yorok zanzanok .
5b. totat jjat quat cat .	11b. wat nnat arrat mat zanzanat .
6a. lalok sprok izok jok stok .	12a. lalok rarok nok izok hihok mok .
6b. wat dat krat quat cat .	12b. wat nnat forat arrat vat gat .

Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan: **farok** **errrok** **hihok** **yorok** **clok** kantok ok-yurp

1a. ok-voon ororok sprok .	7a. lalok farok ororok lalok sprok izok enemok .
1b. at-voon bichat dat .	7b. wat jjat bichat wat dat vat eneat .
2a. ok-drubel ok-voon anak plok sprok .	8a. lalok brok anak plok nok .
2b. at-drubel at-voon pippat rrat dat .	8b. iat lat pippat rrat nnat .
3a. erok sprok izok hihok ghirok .	9a. wiwok nok izok kantok ok-yurp .
3b. totat dat arrat vat hilat .	9b. totat nnat quat oloat at-yurp .
4a. ok-voon anak drok brok jok .	10a. lalok mok nok yorok ghirok clok .
4b. at-voon krat pippat sat lat .	10b. wat nnat gat mat bat hilat .
5a. wiwok farok izok stok .	11a. lalok nok errrok hihok yorok zanzanok .
5b. totat jjat quat cat .	11b. wat nnat arrat mat zanzanat .
6a. lalok sprok izok jok stok .	12a. lalok rarok nok izok hihok mok .
6b. wat dat krat quat cat .	12b. wat nnat forat arrat vat gat .

Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan: **farok** **errok** **hihok** **yorok** **clok** kantok ok-yurp

1a. ok-voon ororok sprok .	7a. lalok farok ororok lalok sprok izok enemok .
1b. at-voon bichat dat .	7b. wat jjat bichat wat dat vat eneat .
2a. ok-drubel ok-voon anak plok sprok .	8a. lalok brok anak plok nok .
2b. at-drubel at-voon pippat rrat dat .	8b. iat lat pippat rrat nnat .
3a. erok sprok izok hihok ghirok .	9a. wiwok nok izok kantok ok-yurp .
3b. totat dat arrat vat hilat .	9b. totat nnat quat oloat at-yurp .
4a. ok-voon anak drok brok jok .	10a. lalok mok nok yorok ghirok clok .
4b. at-voon krat pippat sat lat .	10b. wat nnat gat mat bat hilat .
5a. wiwok farok izok stok .	11a. lalok nok errok hihok yorok zanzanok .
5b. totat jjat quat cat .	11b. wat nnat arrat mat zanzanat .
6a. lalok sprok izok jok stok .	12a. lalok rarok nok izok hihok mok .
6b. wat dat krat quat cat .	12b. wat nnat forat arrat vat gat .

process of elimination

Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan: **farok** **errok** **hihok** **yorok** **clok** kantok ok-yurp

1a. ok-voon ororok sprok .	7a. lalok farok ororok lalok sprok izok enemok .
1b. at-voon bichat dat .	7b. wat jjat bichat wat dat vat eneat .
2a. ok-drubel ok-voon anak plok sprok .	8a. lalok brok anak plok nok .
2b. at-drubel at-voon pippat rrat dat .	8b. iat lat pippat rrat nnat .
3a. erok sprok izok hihok ghirok .	9a. wiwok nok izok kantok ok-yurp .
3b. totat dat arrat vat hilat .	9b. totat nnat quat oloat at-yurp .
4a. ok-voon anak drok brok jok .	10a. lalok mok nok yorok ghirok clok .
4b. at-voon krat pippat sat lat .	10b. wat nnat gat mat bat hilat .
5a. wiwok farok izok stok .	11a. lalok nok errok hihok yorok zanzanok .
5b. totat jjat quat cat .	11b. wat nnat arrat mat zanzanat .
6a. lalok sprok izok jok stok .	12a. lalok rarok nok izok hihok mok .
6b. wat dat krat quat cat .	12b. wat nnat forat arrat vat gat .

cognate?

Centauri/Arcturan [Knight, 1997]

Your assignment, put these words in order: { **jjat**, **arrat**, **mat**, **bat**, **oloot**, **at-yurp** }

1a. ok-voon ororok sprok .	7a. lalok farok ororok lalok sprok izok enemok .
1b. at-voon bichat dat .	7b. wat jjat bichat wat dat vat eneat .
2a. ok-drubel ok-voon anak plok sprok .	8a. lalok brok anak plok nok .
2b. at-drubel at-voon pippat rrat dat .	8b. iat lat pippat rrat nnat .
3a. erok sprok izok hihok ghirok .	9a. wiwok nok izok kantok ok-yurp .
3b. totat dat arrat vat hilat .	9b. totat nnat quat oloat at-yurp .
4a. ok-voon anak drok brok jok .	10a. lalok mok nok yorok ghirok clok .
4b. at-voon krat pippat sat lat .	10b. wat nnat gat mat bat hilat .
5a. wiwok farok izok stok .	11a. lalok nok errok hihok yorok zanzanok .
5b. totat jjat quat cat .	11b. wat nnat arrat mat zanzanat .
6a. lalok sprok izok jok stok .	12a. lalok rarok nok izok hihok mok .
6b. wat dat krat quat cat .	12b. wat nnat forat arrat vat gat .

zero fertility

It's Really Spanish/English

Clients do not sell pharmaceuticals in Europe => Clientes no venden medicinas en Europa

1a. Garcia and associates .	7a. the clients and the associates are enemies .
1b. Garcia y asociados .	7b. los clients y los asociados son enemigos .
2a. Carlos Garcia has three associates .	8a. the company has three groups .
2b. Carlos Garcia tiene tres asociados .	8b. la empresa tiene tres grupos .
3a. his associates are not strong .	9a. its groups are in Europe .
3b. sus asociados no son fuertes .	9b. sus grupos estan en Europa .
4a. Garcia has a company also .	10a. the modern groups sell strong pharmaceuticals .
4b. Garcia tambien tiene una empresa .	10b. los grupos modernos venden medicinas fuertes .
5a. its clients are angry .	11a. the groups do not sell zenzanine .
5b. sus clientes estan enfadados .	11b. los grupos no venden zanzanina .
6a. the associates are also angry .	12a. the small groups are not modern .
6b. los asociados tambien estan enfadados .	12b. los grupos pequenos no son modernos .

Data available

Many languages

- Europarl corpus has all European languages
 - <http://www.statmt.org/europarl/>
 - From a few hundred thousand sentences to a few million
- French/English from French parliamentary proceedings
- Lots of Chinese/English and Arabic/English from government projects/interests
 - Chinese-English: 440 million words (15-20 million sentence pairs)
 - Arabic-English: 790 million words (30-40 million sentence pairs)
- Smaller corpora in many, many other languages

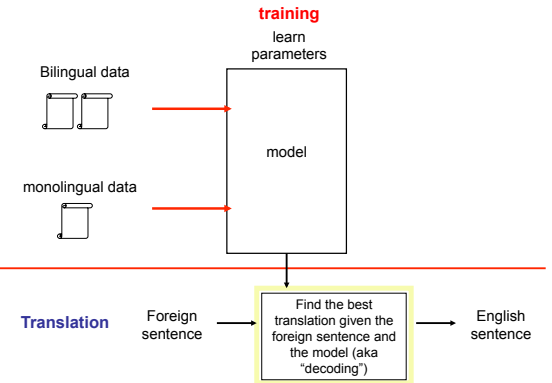
Lots of monolingual data available in many languages

Even less data with multiple translations available

Available in limited domains

- most data is either news or government proceedings
- some other domains recently, like blogs

Statistical MT Overview



Statistical MT

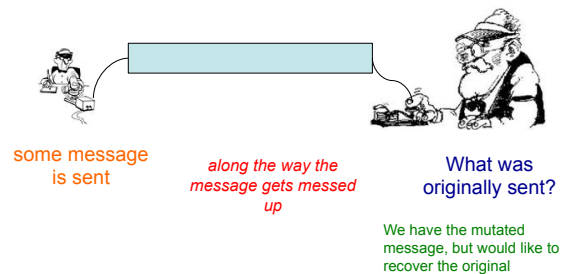
We will model the translation process probabilistically

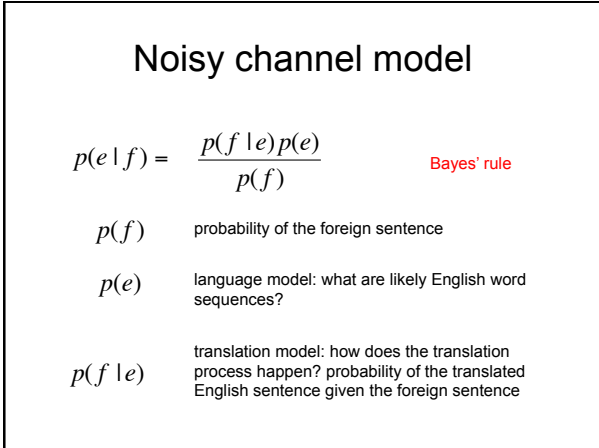
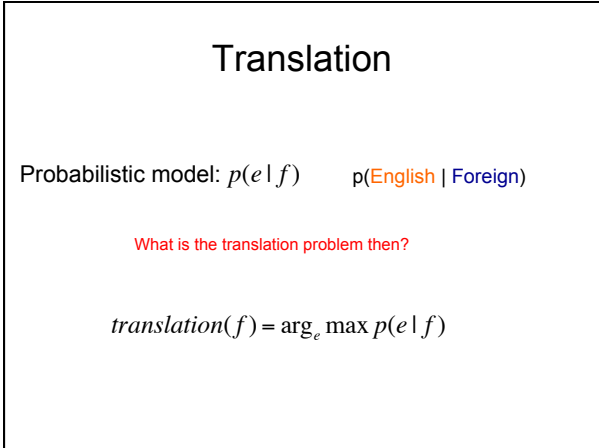
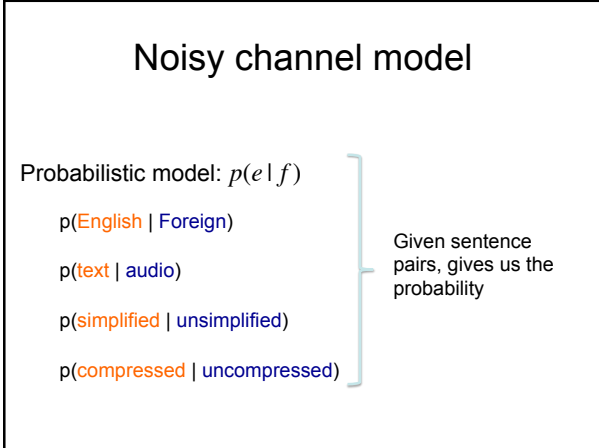
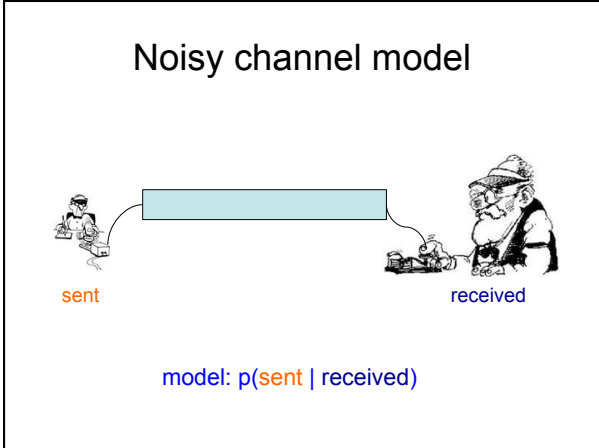
Given a foreign sentence to translate, for any possible English sentence, we want to know the probability that the sentence is a translation of the foreign sentence

If we can find the most probable English sentence, we're done

$$p(\text{english sentence} \mid \text{foreign sentence})$$

Noisy channel model





Noisy channel model

$p(e | f) = p(f | e)p(e)$ Bayes' rule

~~$p(f)$~~ probability of the foreign sentence why?

$p(e)$ language model: what are likely English word sequences?

$p(f | e)$ translation model: how does the translation process happen? probability of the translated English sentence given the foreign sentence

Noisy channel model

$p(e | f) = p(f | e)p(e)$ Bayes' rule

~~$p(f)$~~ probability of the foreign sentence why?

$translation(f) = \arg_e \max \frac{p(f | e)p(e)}{p(f)} = \arg_e \max p(f | e)p(e)$

this is a constant for any given f

Noisy channel model

model $p(e | f) \propto p(f | e)p(e)$

translation model language model

how do English sentences get translated to foreign? what do English sentences look like?

Translation model

The models define probabilities over inputs $p(f | e)$

Morgen fliege ich nach Kanada zur Konferenz

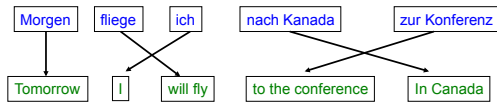
Tomorrow I will fly to the conference in Canada

What is the probability that the English sentence is a translation of the foreign sentence?

Translation model

The models define probabilities over inputs

$$p(f|e)$$



- What is the probability of a foreign word being translated as a particular English word?
- What is the probability of a foreign foreign phrase being translated as a particular English phrase?
- What is the probability of a word/phrase changing ordering?
- What is the probability of a foreign word/phrase disappearing?
- What is the probability of a English word/phrase appearing?

Translation model

The models define probabilities over inputs

$$p(f|e)$$

$$p(\text{Morgen fliege ich nach Kanada zur Konferenz} | \text{Tomorrow I will fly to the conference in Canada}) = 0.1$$

$$p(\text{Morgen fliege ich nach Kanada zur Konferenz} | \text{I like peanut butter and jelly}) = 0.0001$$

Language model

The models define probabilities over inputs

$$p(e)$$

Tomorrow I will fly to the conference in Canada

What is a probability distribution?

A probability distribution defines the probability over a space of possible inputs

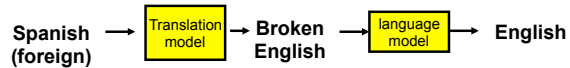
For the language model, what is the space of possible inputs?

- A language model describes the probability over **ALL** possible combinations of English words

For the translation model, what is the space of possible inputs?

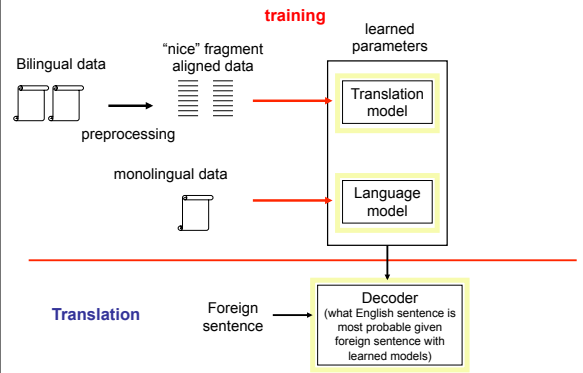
- **ALL** possible combinations of foreign words with **ALL** possible combinations of English words

One way to think about it...



Que hambre tengo yo → What hunger have I,
 Hungry I am so, → I am so hungry
 I am so hungry,
 Have I that hunger ...

Statistical MT Overview



Basic Model, Revisited

$$\operatorname{argmax}_e P(e | f) =$$

$$\operatorname{argmax}_e P(e) \times P(f | e) / P(f) =$$

$$\operatorname{argmax}_e P(e) \times P(f | e)$$

Basic Model, Revisited

$$\operatorname{argmax}_e P(e | f) =$$

$$\operatorname{argmax}_e P(e) \times P(f | e) / P(f) =$$

$$\operatorname{argmax}_e P(e)^{2.4} \times P(f | e) \quad \dots \text{works better!}$$

Basic Model, Revisited

$$\operatorname{argmax}_e P(e | f) =$$

$$\operatorname{argmax}_e P(e) \times P(f | e) / P(f)$$

$$\operatorname{argmax}_e P(e)^{2.4} \times P(f | e) \times \text{length}(e)^{1.1}$$

\

Rewards longer hypotheses, since these are unfairly punished by $P(e)$

Basic Model, Revisited

$$\operatorname{argmax}_e P(e)^{2.4} \times P(f | e) \times \text{length}(e)^{1.1} \times \text{KS}^{3.7} \dots$$

Lots of knowledge sources vote on any given hypothesis.

"Knowledge source" = "feature function" = "score component".

A feature function simply scores a hypothesis with a real value.

(May be binary, as in "e has a verb").

Problems for Statistical MT

Preprocessing

- How do we get aligned bilingual text?
- Tokenization
- Segmentation (document, sentence, word)

Language modeling

- Given an English string e , assigns $P(e)$ by formula

Translation modeling

- Given a pair of strings $\langle f, e \rangle$, assigns $P(f | e)$ by formula

Decoding

- Given a language model, a translation model, and a new sentence f ... find translation e maximizing $P(e) \times P(f | e)$

Parameter optimization

- Given a model with multiple feature functions, how are they related? What are the optimal parameters?

Evaluation

- How well is a system doing? How can we compare two systems?