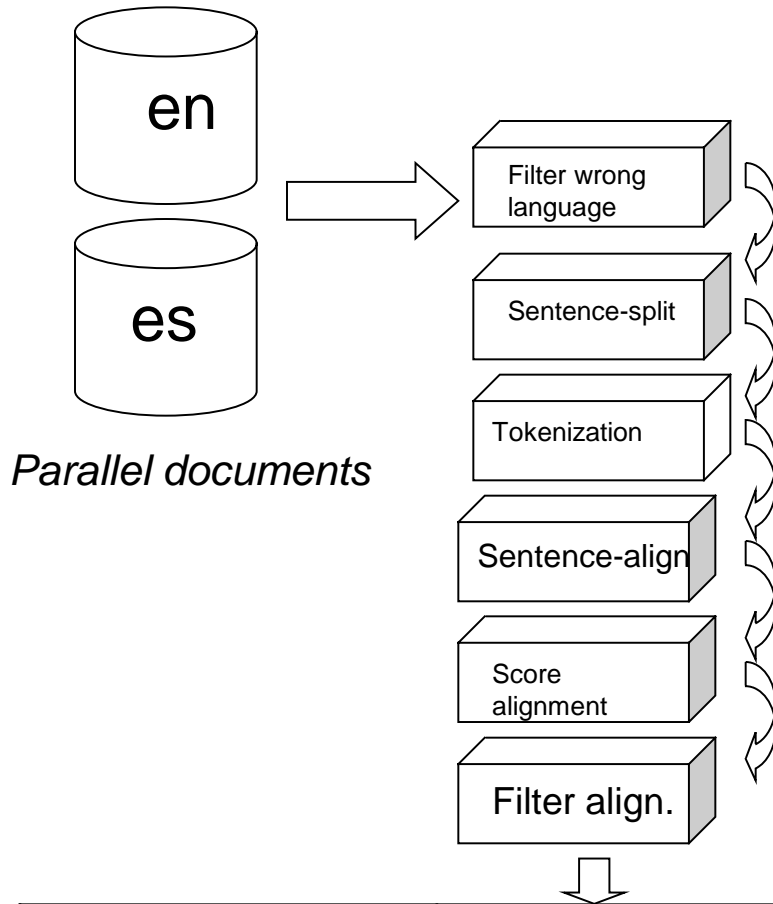


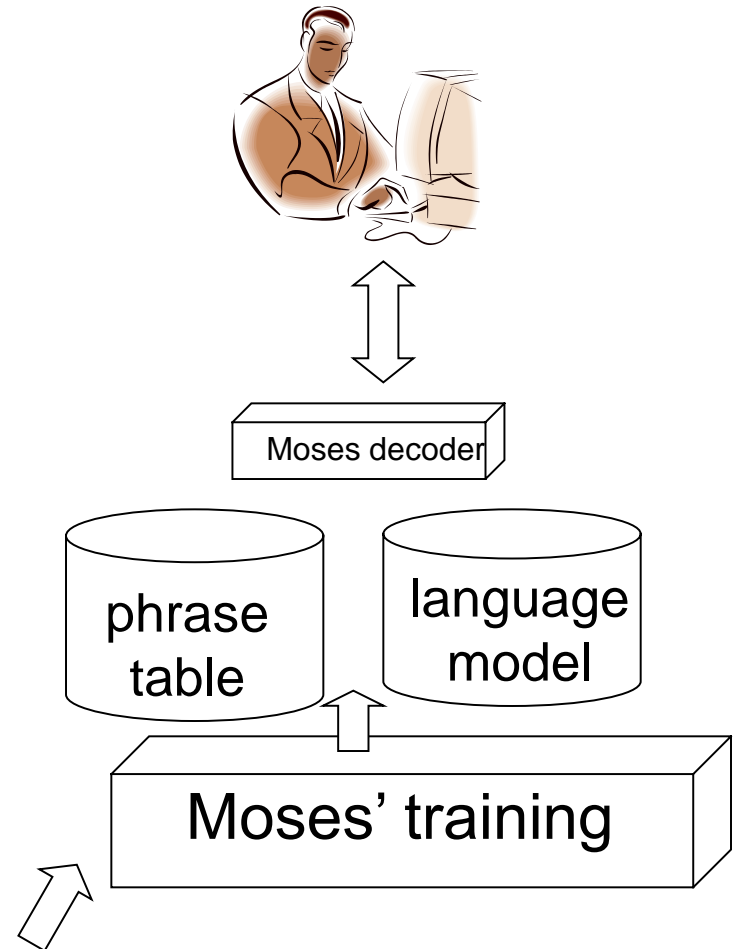
# Lecture 5.4: MT preparation steps

How to prepare your texts for MT,  
pre and post processing

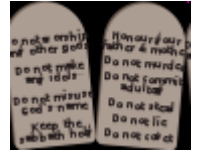
# SMT lifecycle



En	Es
Strengthening of forum for human dignity : legal aid	Fortalecimiento del foro para la dignidad humana - asistencia jurídica
must respect all aspects of human dignity	debe respetar todos los aspectos de la dignidad humana
should fully respect human dignity	se deben respetar plenamente la dignidad humana



# Moses input format



- 2 text files (utf-8), called bitext
- Each line is a tokenized sentence (hopefully short <80w)

Each line in first file has to be the translation of corresponding line in second file

- No id, no classification etc.

# Tokenization



- Moses translates tokens (usually words separated by space)  
Tony ≠ Tony's ≠ Tony. ≠ tony ≠ Tony, ≠ Tony's  
Japanese/Chinese tokenization needed

# Tokenizer, options



For example, our tokenizer in WIPO does:

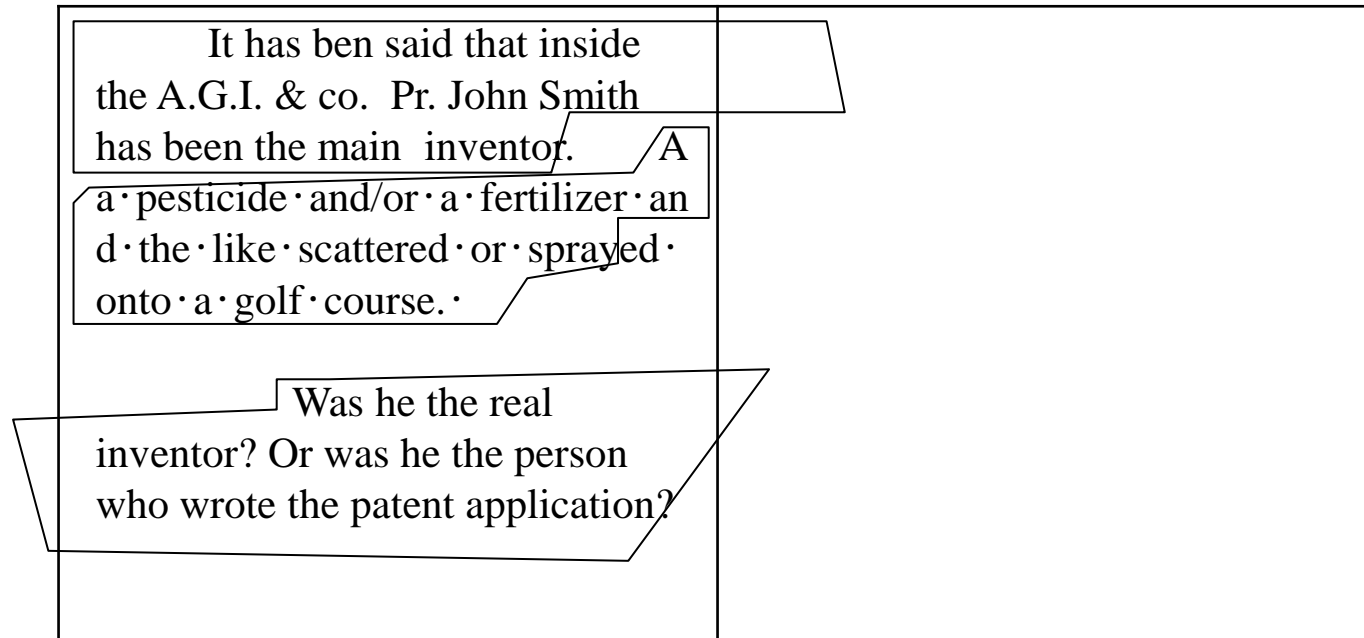
- Extract email address as one single token
- Extract full URLs as one single token (tagged as HOST)
- Recognizes Xml/Html tags as one token
- Reduces acronyms without dots (“U.N.” become “UN”)
- French apostrophe: “jusqu’alors” becomes “jusqu’\_alors”
- English possessive case: “UN’s opinion” becomes “UN\_’s\_opinion”, (but “al-Ma’sara” becomes “al\_-\_ma’sara”)
- Special recasing: “The New York UN HeadQuarter” → “the New York UN headquarter”
- Groups references to figures as one token (eg. “(1)” not “(\_1\_)”)
- Recognize "and / or" as one single token
- Special case of greek letters "α- and/or β-amino acid" becomes "alpha - and/or beta - amino acid “
- Harmonize language specific diacritics? Latin “ae” => ae, delete Arabic short vowels?

# Tokenizer(2)

- Chinese tokenizer: smartcn (small adaptations)
  - “演奏音乐的娱乐装置” => “演奏 / 音乐 / 的 / 娱乐 / 装置”
- German decompounding
  - “Lebenshaltungskostenausgleich” ➔ “Lebens- haltungs- kosten- ausgleich”
- Arabic prefix splitting
  - splits prefixes "ال", "وال", "لل", "بال", "كال"
- Indic languages tokenization
  - [http://anoopkunchukuttan.github.io/indic\\_nlp\\_library/](http://anoopkunchukuttan.github.io/indic_nlp_library/) python script
  - Be careful: the Moses' tokenizer does not handle properly Indic languages
  - For the labs, I quickly adapted the Moses tokenizer (Perl script), see [icon2016/labs/tokenizer.pl](http://icon2016/labs/tokenizer.pl)

# Sentence splitting

- MT works with sentences, a document must be split in sentences



# Sentence splitting

- Usually using full stop as a first approximation
  - However, abbreviations may also contains full stops.
  - E.g.: “As Pr. Smith, no. 123, 13 dec. 2015”
- Depends on the language
  - Chinese “。”, Devanagari “।”, Armenian “:”, Amharic “፥” ...
  - In Thai, a space is the equivalent of a full stop



# Sentence alignment

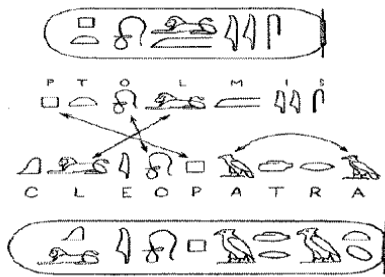
- Documents may contain different sentence structure (sentences are not mapped one-to-one)
- Sentence alignment tools are needed
  - Count n. of words, find best 1-0,0-1,1-1,1-2,2-1 combinations to align [Brown et al.1991])
  - Use bilingual dictionary
  - Use cognates, names, numbers
  - See [Singh & Hussain 2005]
- More on <http://www.statmt.org/survey/Topic/SentenceAlignment>

# Sentence alignment example

<p>a · pesticide · and/or · a · fertili zer · and · the · like · scattered · or · spr ayed · onto · a · golf · course · (1) · , · to gether · with · most · of · the · sprayed · water · or · rain · water · , · are · permated · into · a · turf · (15A) · , · a · thick · sand · layer · (15) · having · high · water · pe rmeability · , · whereby · only · the · rain · water · filtered · by · the · golf · course · is · discharged · to · the · outside</p>	<p>un · pesticide · et/ou · un · fertilisant · · , · ou · un · produit · similaire · , · dispersé · ou · pulvérisé · sur · un · parcours · de · golf · (1) · , · avec · la · majeure · partie · de · l' · eau · pulvérisée · ou · de · l' · eau · de · pluie · , · qui · sont · imprégnés · dans · un · gazon · (15A) · , · une · épaisse couche · de · sable · (15) · ayant · une · forte · perméabilité · à · l' · eau · , · Selon · ce · procédé · seule · l' · eau · de · pluie · filtrée · par · le · sol · du · parcours · de · golf · se · déverse · à · l' · extérieur</p>
--	--

# Document alignment

- Challenges with getting parallel sentences from comparable documents



[0009] Zur Lösung dieses Problems sieht die Erfindung vor, dass bei dem eingangs genannten Tonträger wenigstens eine der Justierdurchbrechungen auf der Oberseite des Tonträgers von einem nach oben aus der Oberseite vorstehenden Wulstring, der z. B. eingepreßt oder

surround it a small distance of about 1.0 to 1.5 mm, especially 1.2 mm, thus forming a flat shoulder which contributesto defining the perforation	1-1	99.28%	einen geringen Abstand von etwa 1,0 bis 1,5 mm, insbesondere 1,2 mm, und bildet so einen flachen Absatz, was zur Präzisierung der Durchbrechung beiträgt
the annular bead should have a height of 0.2 to 0.3 mm, especially of 0.2 mm, and a width of about 1-1.5 mm, especially 1.2 mm. it is advisable to provide the two perforations in the lower edge of the sound carrier with such an annular bead	1-2	48.85%	Der Wulstring sollte eine Höhe von 0,2 bis 0,3 mm, insbesondere von 0,2 mm, und eine Breite von ca. 1 - 1,5 mm, insbesondere 1,2 mm haben. Zweckmässigerweise werden zwei Durchbrüche am unteren Rand des Tonträgers mit diesem Wulstring ausgestattet
such an annular bead does not only make the perforations better visible, they also render the alignment of the feet of the sound reproducing device easier because they practically fall into the larger annular bead and slip into the precise adjustment perforations automatically	1-1	90.33%	Durch diesen Wulstring sind die Durchbrüche nicht nur besser sichtbar, sie erleichtern auch das Justieren der Füße des Tonwiedergabegerätes, da die Justierfüsse in dem grösseren Wulstring praktisch hineinfließen und zwangsweise in die präzisen Justierdurchbrechungen schlüpfen
such a perforation surrounded by an annular bead can actually be compared to a funnel	1-1	21.93%	Man kann diese von einem Wulstring umgebenen Durchbrüche praktisch mit einem Trichter vergleichen.
of course, all the perforations in a sound carrier can be provided with beads	1-1	30.73%	Selbstverständlich können auch alle Durchbrechungen eines Tonträgers mit Wulsten versehen werden
they also prevent an unintentional dislocation of the sound reproducing device during possibly desired interruptions of a reproduction while studying educational information. This is important as the sound reproducing device has to be actuated frequently also during reproduction, for example for interruptions or for a repeated reproduction which naturally involves the danger of dislocation	1-2	50.66%	Diese Wülste verhindern auch das unbeabsichtigte Verrücken des Tonwiedergabegerätes bei evtl. gewollten Unterbrechungen der Wiedergabe beim Studium von Lerntexten. Das ist wichtig, weil das Tonwiedergabegerät auch während der Wiedergabe häufig, beispielsweise für Unterbrechungen oder Wiedergabewiederholungen zu betätigen ist und dabei an sich die Gefahr des Verrückens besteht.
as an alternative to this solution, the two lower adjustment perforations may be widened and elongated in the direction away from the other adjustment perforations toward the outside or in the direction of the same toward the inside and surround the respective aligning foot to be inserted into them with large clearance by means of the wide region created in this manner, so that said aligning foot is movable to abut the narrow region and the other aligning foot / feet is / are in a position to engage the other adjustment perforation / perforations in this position	1-1	100%	Alternativ zu dieser Lösung kann auch ein Teil der Justierdurchbrechungen, z. B. zwei untere Justierdurchbrechungen, in Richtung von den anderen Justierdurchbrechungen nach aussen weg oder zu diesen nach innen hin verbreitert und verlängert ausgebildet sein und den in sie einzuführenden Justierfuss mit dem so geschaffenen breiten Bereich mit grossem Spiel so umschliessen, dass dieser Justierfuss zur Anlage an den engen Bereich verschieblich ist und die anderen Justierfüsse in dieser Lage in die anderen Justierdurchbrechungen, die die Justierfüsse im engsten Bereich mit geringem Spiel umschliessen, einzugreifen vermögen
a plane guiding funnel is formed in this manner which passes the sound reproducing device placed in position into the round perforations of the sound carrier	1-1	99.86%	So ist ein ebener Führungstrichter gebildet, der das aufgesetzte Tonwiedergabegerät in die runden Durchbrechungen des Tonträgers leitet
	0-1		Zusammen mit den dann zum Eingriff kommenden anderen Justierfüssen ist das Tonwiedergabegerät dann lagestabil gehalten.
	0-1		Ein älterer bekannter Vorschlag (US Patentschrift 4,298,967) ist von der Praxis als ungeeignet verworfen worden, der im Mittelpunkt der Tonrille eines folienartigen, dünnen Tonträgers eine schuhartige zur Mitte und nach unten sich verengende Vertiefung vorsah, in die ein zentraler Justierstift des Tonwiedergabegerätes bis zum Grund in die Justierstellung geführt wird, in diese aber wegen Fehlen der Arretierungsmöglichkeit nicht lagestabil festgehalten ist
	0-1		Ausserdem verdeckt das Tonwiedergabegerät jegliche Sicht zur schuhartigen Vertiefung.
the walls of the enlarged adjustment perforations should have the same configuration on the side facing toward or away from the other adjustment perforations as the aligning feet to be inserted	1-1	57.12%	Die Wände der vergrösserten Justierdurchbrechungen sollten auf der den anderen Justierdurchbrechungen zugewandten oder abgewandten Seite die gleiche Gestalt haben wie die einzusetzenden Justierfüsse
as a result, when centered, the feet are in each case in the exactly defined			Dadurch befinden sich die Füße lagestabil in der zentrierten Lage jeweils an den genau

ann sich unmittelbar an die Durchbrechung d von etwa 1,0 bis 1,5 mm, insbesondere 1,2 mm, Präzisierung der Durchbrechung beiträgt. Der insbesondere von 0,2 mm, und eine Breite von brüche am unteren Rand des Tonträgers mit ulstring sind die Durchbrüche nicht nur besser üsse des Tonwiedergabegerätes, da die hineinfließen und zwangsweise in die präzisen ese von einem Wulstring umgeben en stehen.

chbrechungen eines Tonträgers mit Wulsten n das unbeabsichtigte Verrücken des brechungen der Wiedergabe beim Studium von aberät auch während der Wiedergabe häufig, rgabewiederholungen zu betätigen ist und

n Teil der Justierdurchbrechungen, z. B. zwei den anderen Justierdurchbrechungen nach ert und verlängert ausgebildet sein und den in affenen breiten Bereich mit grossem Spiel so n den engen Bereich verschieblich ist und die n Justierdurchbrechungen, die die Justierfüsse n, einzugreifen vermögen. So ist ein ebener onwiedergabegerät in die runden en mit den dann zum Eingriff kommenden it dann lagestabil gehalten.

entschrift 4,298,967) ist von der Praxis als tkt der Tonrille eines folienartigen, dünnen unten sich verengende Vertiefung vorsah, in die bis zum Grund in die Justierstellung geführt gsmöglichkeit nicht lagestabil festgehalten ist. icht) Sicht zur schuhartigen Vertiefung.

hbrechungen sollten auf der den anderen ewandten Seite die gleiche Gestalt haben wie n sich die Füße lagestabil in der zentrierten n, an denen sie sich auch dann befinden, wenn ären. Die nicht vergrößerten schuhartiger, also einen von etwa 5 - 6 mm en etwa zu knapp der Hälfte die gleiche Gestalt, eine größte Querabmessung von etwa 9 - 10 te des Tonträgers offen. An sich genügen für den eine rund und die andere trichterförmig chbrechungen vorgesehen, von denen dann zwei

# Text specificities

- Challenges with parallel sentences from documents (here lists sorted alphabetically)

China	
Colombia	
Comoros (OA) <sup>2,3</sup>	
Congo (OA) <sup>2</sup>	
Costa Rica	
Côte d'Ivoire (OA) <sup>2</sup>	
Croatia (EP)	
Cuba	
Cyprus (EP) <sup>2</sup>	
Czech Republic (EP)	
Democratic People's Republic of Korea	
Denmark (EP)	

Chine	
Chypre (EP) <sup>2</sup>	
Colombie	
Comores (OA) <sup>2,3</sup>	
Congo (OA) <sup>2</sup>	
Costa Rica	
Côte d'Ivoire (OA) <sup>2</sup>	
Croatie (EP)	
Cuba	
Danemark (EP)	
Dominique	

# Clean long lines



- Filter out sentences having more than 80 words
- Ratio (# src words/ # dst words) more than 9

```
perl ~/mosedecoder/script/bin/clean-corpus-n.perl training en es  
cleaned 1 80
```

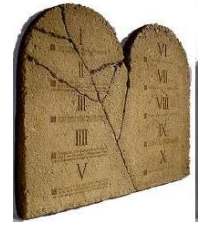
Word alignment (giza) complexity is linear according to the number of sentences, however it is quadratic according to the number of tokens per sentence

# Trainset / devset / testset



- Trainset: Set of sentences used by your machine to learn (& language model generation)
- Devset: Set of sentences used to optimize the model parameters (Moses “MERT”)
- Testset: Set of sentences used for comparison between automatic translation and human translation (e.g. BLEU)

# Trainset / devset / testset few commandments



- Never mix trainset and testset  
(Even for the language model)
- Always test on held-out documents
- Try to have a representative test/dev set  
(In some context: translate new documents)
  - Test on newest documents
  - Develop on newest documents
- Use appropriate evaluation measure (human evaluation, BLEU etc.), see lecture 6