

# Ontology as Interlingua

# Crosslinguistic WordNets

Starting in late 1990s, WordNets were built for languages other than English

Genetically and typologically unrelated languages:  
Turkish, Hindi, Chinese, Korean, Basque, Xhosa, Arabic, Latin... (currently >70)

<http://www.globalwordnet.org>

# Wordnets

Entire families of wordnets: EuroWordNet,  
BalkaNet, IndoWordNet, AsianWordNet,  
AfricanWordNet,...

Local product: MultiWordNet

# Wordnets in the world

## Motivations:

- Natural Language Processing applications that require word sense discrimination and disambiguation within and across languages
- Crosslingual comparison of lexical categories
- Interesting by-product: discover language-specific lexical gaps

# What is universal?

- Surely not all “concepts”:
- English has many verbs of walking (*slouch, strut, stroll, amble, prance, sneak, march,..*) and walking/running (*hop, skip, bounce,..*)
- No 1:1 crosslingual encoding of concepts
- But is the network structure universal? Can all words in all languages be connected?
- Are the relations universal (if so, this would strengthen their cognitive reality)

# Classes

- Crosslingual construction reveals potentially meaningful classes
- Classifiers (Chinese, Bantu)
- Represent roots as underspecified semantic categories separately from words in semitic languages

# Crosslinguistic WordNets

Some are manually constructed

--independently from PWN, mapped later (“Merge” method)

or

--translated directly from PWN (“Expand” method)

First method is considered easier, more accurate (why?)

Other wordnets are constructed semi-automatically

# Crosslinguistic Wordnets

All new wordnets are mapped to the Princeton WordNet , which serves as a general model and as the link for connecting the wordnets



# Mapping words and synsets across multilingual WordNets

First set of eight foreign-language WNs (EuroWordNet; Vossen 1998) were built with reference to Princeton WordNet

Princeton WN as the hub (“interlingual index”)

Each synset in each WN was linked to a “record” (PWN synset identifier) in the index

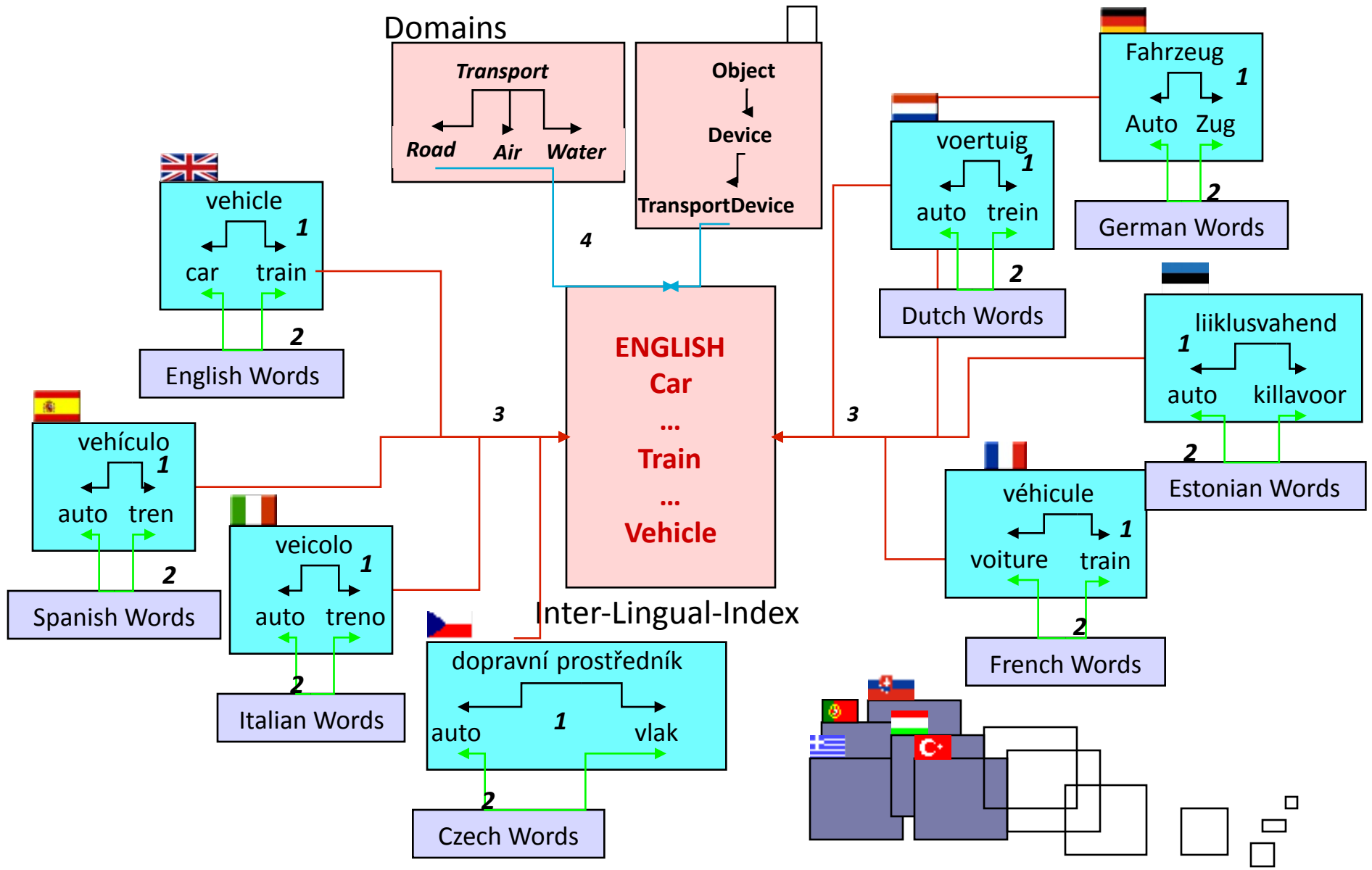
Crosslingual mapping of words and synsets proceeds via the index

# Mapping words and synsets across multilingual WordNets

The Interlingual Index is a flat, unstructured list

Princeton WN's **structure** is not imported

Only the language-specific wordnets have relations and form networks



# Mismatches in multilingual WordNets

Concepts not lexicalized in English required the creation of new records in the ILI (w/out English synsets or synsets in some other wordnets)

E.g., Arabic lexically distinguishes more kinds of *cousin* than English; thus the ILI needs appropriate placeholders (records)

Xhosa time expressions:

*the\_time\_of\_day\_when\_you\_are\_beautiful*

*the\_time\_of\_day\_when\_you\_see\_the\_horns\_of\_the  
\_cattle\_against-the\_sky*

# Mismatches in multilingual WordNets

Conversely, some languages lack the equivalents of English words:

--Dutch lacks *container* but lexicalizes kinds (hyponyms) of *container* (*box, bag, bottle, bowl...*)

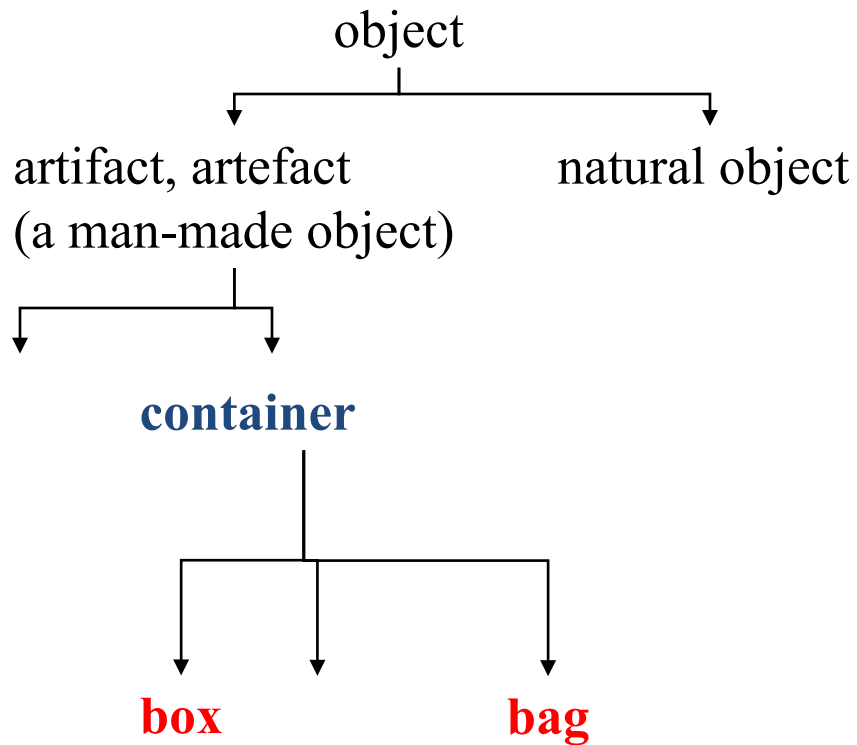
Respective hierarchies reflects this difference; Dutch wordnet “skips” a level

Du. *bag, box...* => *artifact*

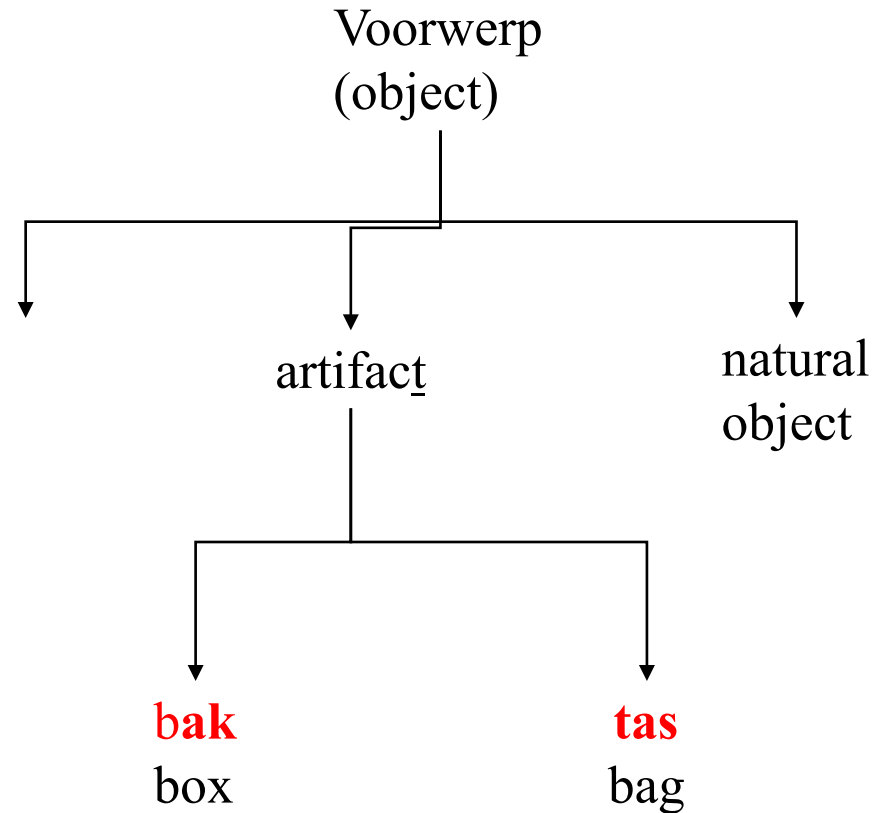
Engl. *bag, box...* => *container* => *artifact*

# English-Dutch snippet

## English Wordnet



## Dutch Wordnet



# Multilingual WordNets

Interlingual Index in EuroWordNet is biased towards English

Could skew coverage of new wordnets, esp. those that are translated from English

Some mapped synsets aren't really equivalent

# Interlingua

Solution: replace index by *language-independent, formal* ontology that can accommodate the lexicons of all languages

Meanings are stated as axioms in logical form

Axioms are machine-readable

Interlingual ontology enables automatic reasoning and inferencing, within and across languages



# Some proposed ontologies

- SUMO
- DOLCE
- KYOTO