

## Setting



- similar languages (e.g. Slovak and Czech)
- no parallel data, no supervised source data

## Goal



- simulate human-like crosslingual understanding
- perform machine translation without parallel data

## Reasoning



- there are ~7000 languages in the world
- required data easy to get for ~100 languages

## Method

**Idea:** translate each source word by the **most similar** target word

**Assumptions:** corresponding words:

- are string-wise similar
- have similar length
- have similar frequency

$$\text{sim}(s, t) = \text{sim}_{JW}(s, t) \cdot \text{sim}_{len}(s, t) \cdot \text{sim}_{freq}(s, t)$$

### String-wise similarity

- based on Jaro-Winkler edit distance
- prefixes more important
- also computed on ASCII transliterations
- also computed on devowelled words (strip AEIOUY)

$$\text{sim}_{JW}(s, t) = \text{sim}_{jw}(s, t) \cdot \text{sim}_{jwT}(s, t) \cdot \text{sim}_{jwD}(s, t) \cdot \text{sim}_{jwDT}(s, t)$$

### Length similarity

- also computed on devowelled words

$$\text{sim}_{len}(s, t) = \frac{1}{1 + |\text{len}(s) - \text{len}(t)|}$$

### Frequency similarity

- computed on available monolingual corpora
- normalization to account for differing corpora sizes

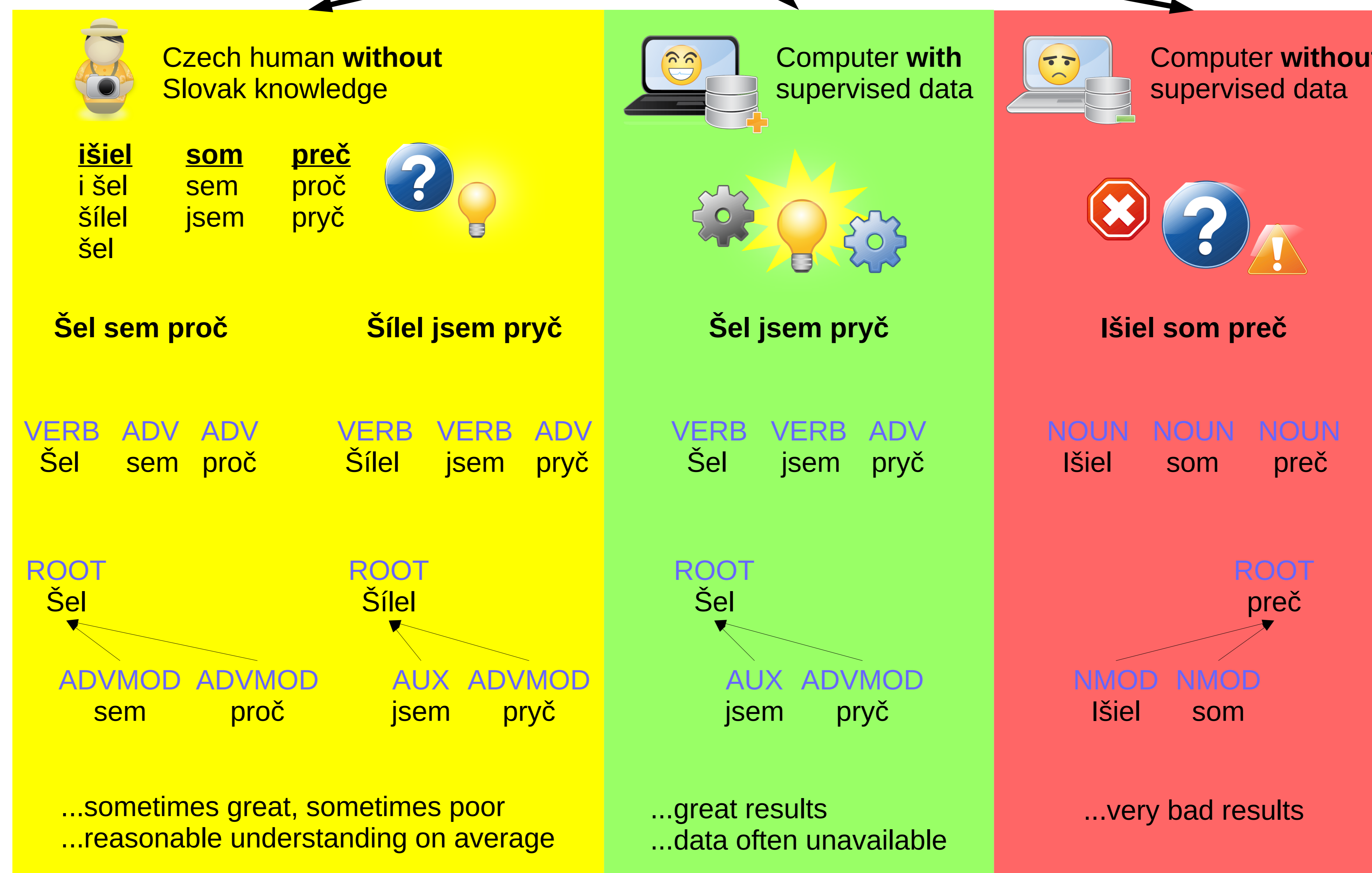
$$\text{sim}_{freq}(s, t) = \frac{1}{1 + |\log(\text{freq}(s)) - \log(\text{freq}(t))|}$$

### Efficiency

- cannot go over all target words for each source word
- leads to e.g. 1 word/minute
- goal: e.g. 1 word/second
- early stopping
- start with frequent candidates
- stop once  $\text{sim}_{freq}$  is lower than current highest sim
- huge speedup
- limited capabilities: hard to add a language model
- word list partitioning
- source and target word must share first 2 consonants
- huge speedup
- introduces search errors, many translations not found
- e.g. "som" cannot be translated to "jsem"

## Motivation

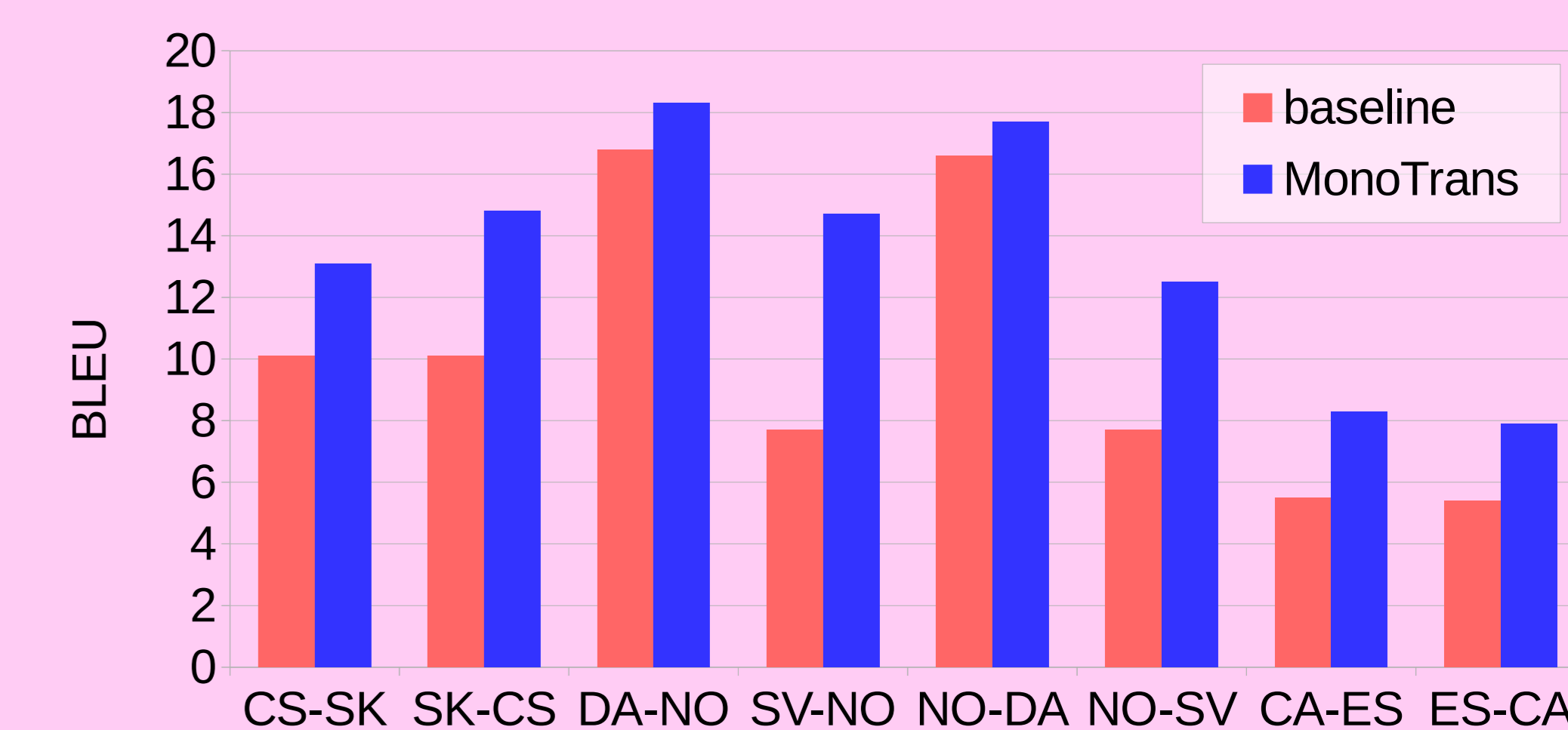
Input: Išiel som preč



## Evaluation

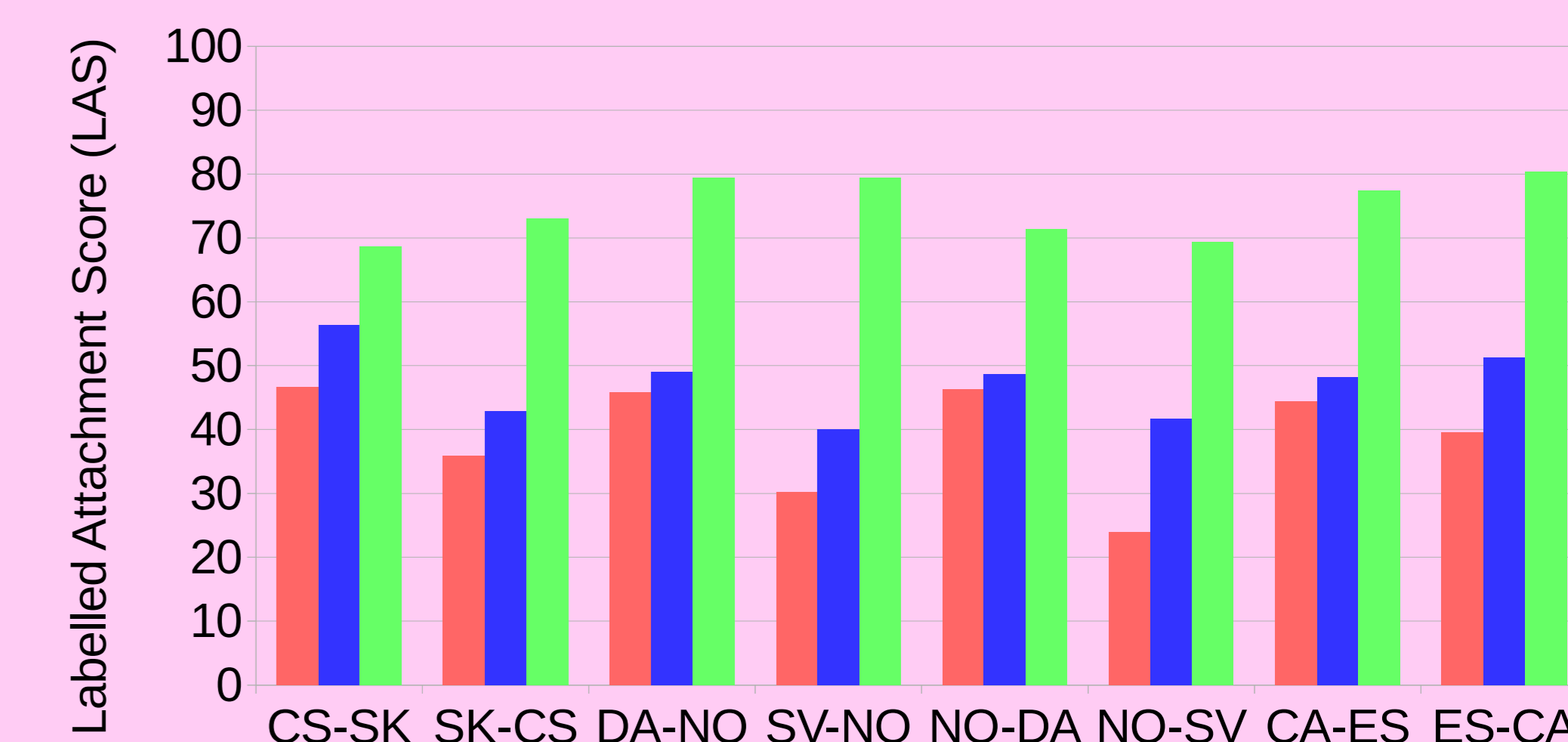
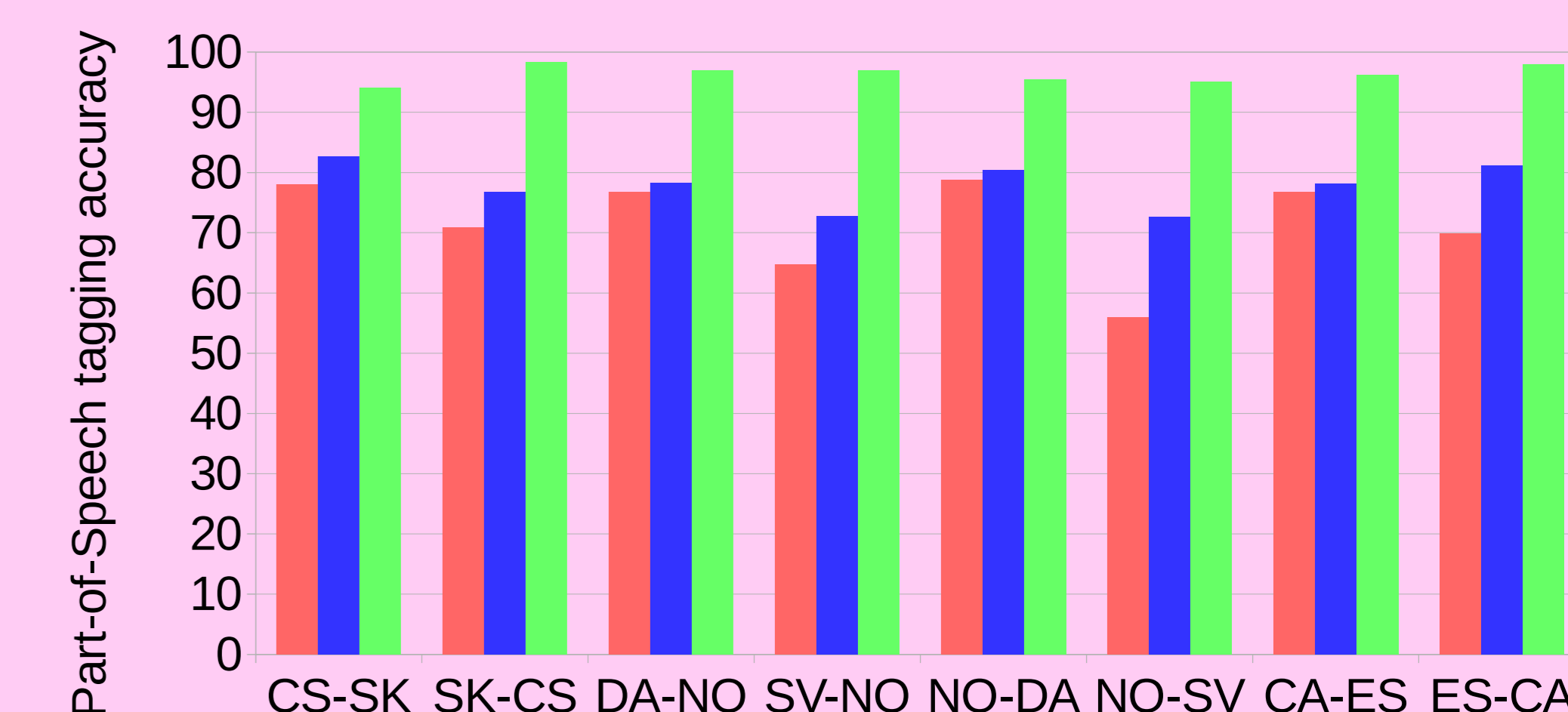
### Intrinsic evaluation: translation quality

- data source: OpenSubtitles2016
- Czech (CS), Slovak (SK)
- Danish (DA), Norwegian (NO), Swedish (SV)
- Catalan (CA), Spanish (ES)
- data size: 10,000 sentences
- baseline: keep source text untranslated
- result: average **43% improvement** over baseline



### Extrinsic evaluation: tagging and parsing

- treebank data: Universal Dependencies v1.4
- translate source UD treebank into target language
- train UDPipe tagger and parser
- evaluate on target UD treebank
- baseline: keep source TB untranslated
- supervised: use target TB for training
- result: average **23% error reduction**



## Sample outputs (SK-CS)

SK: Všetci sme počuli o tom, čo si urobil pre toho chlapca a je mi ľúto, že to nezvládol. CS: Všech sem poučil o tom, čo si urobil pro toho chlapce a je mi líto, že to nezvládol.

SK: Navrhujem , aby si sa zase sústredil na hru. CS: Navrhují, aby si se zase sústredil na hru.

SK: Volá sa to vernosť. CS: Volá se to věrnost.

SK: Keď ma zrazil opitý vodič, ktorého uznali nevinným, bol to Doug, kto mi pomohol zaplatiť za štúdium práva. CS: Kde má zrazil opitý vodil, kterého uznali nevinnou, byl to Doug, který mi pomohl zaplatit za štúdium práva.

SK: Posledné, čo sa mi chce, je strčiť ti loptu do ruky. SK: Zoznámim vás s priateľmi. SK: Počuj, asi som zaujatý, no toto je dokonalá svadba. CS: Poslední, čo se mi chce, je strčte ti lopaty do ruky. CS: Zazvoním vás s parazitem. CS: Počkej, asi sem zajatí, no toto je dokonalý svatba.

SK: Môžem úprimne povedať, že som nikdy nemal priateľa, s ktorým by som sa cítil ako s tebou. SK: Nie! Žiadne arašidové maslo! CS: Můžeme promíne povídat, že sem nikdy neměl promiňte, s kterým by sem se cítil ako s tebou. CS: Ne! Žádné arašidové myslí!

SK: My sa môžeme povalovať sa na pohovkách. SK: Na našom treťom rande povedala, že si chce založiť rodinu, že je pripravená mať deti. CS: My se můžeme považovat se na pohádkách. CS: Na našim treťom rande povídal, že si chce založili rodina, že je připravená mít děti.

SK: Alebo povedať: Doug, počme na pivo. Alebo si pozrime zápas. SK: Mohol by sa aj zvyšok svadobčanov pridať k šťastnému páru? CS: Alebo povídat: Doug, podle na pivo. Alebo si pozřeme zápas. CS: Mohl by se já zvykem svadobčanov prodat k strašnému páru?

SK: Historická pánska jazda s mojím najbližším priateľom. SK: Nekecaj. Myslíš, že ja to nechcem? SK: Vyhlasujem vás za manželov. CS: Historický pánské jízde s mojím nejbližším parazitem. CS: Nekecej. Myslíš, že já to nechceš? CS: Vyhlašujeme vás za manželku.

SK: Snažím sa bojovať so svojimi zbraňami každý deň, každý zápas, po celú dobu svojej kariéry. SK: Nemal by som jej vravieť pravdu? CS: Snažím se bojovat s svojich zbraněmi každý den, každý zápas, po celé dobu svoje kariéry. CS: Neměl by sem její vrahovi pravdu?