

# Posudek vedoucího diplomové práce

Matematicko-fyzikální fakulta Univerzity Karlovy

**Autor práce** Václav Ryšlink

**Název práce** Methods of Input Segmentation for Simultaneous Speech Translation

**Rok odevzdání** 2022

**Studijní program** Informatika **Studijní obor** Umělá inteligence

**Autor posudku** doc. RNDr. Ondřej Bojar, Ph.D. **Role** vedoucí

**Pracoviště** Ústav formální a aplikované lingvistiky

## Text posudku:

Diplomová práce Václava Ryšlinka se zabývá otázkou hledání hranic v mluveném projevu, na kterých je možné přeložit celou úvodní část bez znalosti zbytku. Určení takových hranic má velký význam při vývoji systémů simultánního strojového překladu, tj. překladu v době, kdy ještě nebyl získán celý vstup. Simultánní tlumočníci rozhodnutí, zda překládat nebo ještě čekat, řeší rutinně. Kromě snahy o co nejmenší zpoždění a co nejméně nutných korekcí však musí hledět i na další kritéria, např. své momentální vytížení a kapacitu své okamžité paměti. Studium strojových metod je proto zajímavé i pro srovnání s lidmi.

Práce je psána výbornou angličtinou se zanedbatelným množstvím chyb nebo překlepů. Práce je také nadprůměrně obsáhlá a to i přes kompaktní styl. Text sestává z úvodu, šesti vlastních kapitol a závěru; v přílohách najdeme příklady dělení vstupu od anotátorů, pomocí jedné z navržených metod a grafy k podrobné analýze segmentačních metod.

První kapitola přináší potřebnou jazykovědnou teorii, zejména vypichuje různorodost definic tzv. nejmenších překladových jednotek (minimal translation units). Druhá kapitola představuje oblast simultánního strojového překladu, včetně zásadních otázek vyhodnocování jeho kvality. Třetí kapitola se věnuje podrobnému modelování ideální situace: jak by vstup dělili lidé při úplné znalosti. Součástí této kapitoly je i výborně prodiskutovaný návrh, co a jak anotovat, a souhrnně je představen získaný anotovaný korpus; jde o hodnotný výsledek pro obor počítačové lingvistiky a snad i překladatelsví a tlumočnictví.

Čtvrtá, pátá a šestá kapitola jsou těžištěm Václavovy informatické práce. Čtvrtá kapitola přináší návrh sedmi metod segmentace vstupu vycházejících z různých typů vstupní informace, včetně ručního dělení získaného ve třetí kapitole. Pátá kapitola se věnuje výsledkům dvojího druhu: (1) napřed analyzuje získaná data od anotátorů, tj. výsledky anotace z kapitoly 3, a následně výsledky navržených metod automatického dělení, tj. výsledky metod z kapitoly 4, včetně podrobného zavedení metod hodnocení různých způsobů dělení vstupu. Velmi podnětná je šestá

kapitola, která v diskusi shrnuje zásadní pozorování, rozhodnutí i omezení v průběhu celé práce, a současně otevírá řadu otázek pro příští výzkum jak v oblasti traslatologie a tlumočnictví, tak v oblasti automatického zpracování přirozeného jazyka, a snaží se tyto obory v úloze simultánního překladu těsněji provázat.

Velmi oceňuji Václavovu samostatnost ve všech aspektech práce: při zpracování teorie (zejm. pak z oblasti překladatelství a tlumočnictví, tj. domény aplikace jeho metod, ale i z oblasti strojového překladu a jeho vyhodnocování, což je jen jedna z konkrétních oblastí oboru umělé inteligence), při návrhu anotací pro sběr ideální segmentace, při dohledu nad anotátory, při analýze jejich výstupů, při trénování modelů strojového překladu, při návrhu automatických metod segmentace a zejména pak při analýze jejich chování a následné diskusi. Výsledky práce pokládám za velmi zajímavé a určitě se je následně pokusíme publikovat na vhodné konferenci.

S Václavovou prací jsem celkově mimořádně spokojen a doporučuji, aby byla přijata.

**Práci doporučuji k obhajobě.**

**Práci navrhuji na zvláštní ocenění.**

Navrhuji diplomovou práci Václava Ryšlinka ocenit vhodnou cenou. Práce je zpracována mimořádně pečlivě, i při úsporném psaní je nadprůměrně obsáhlá a Václav samostatně provedl vynikající analýzy. Výsledky z práce jsme dosud nepublikovali, ale určitě o to budeme usilovat, neboť práce výborným způsobem připravuje půdu pro interdisciplinární výzkum v oblasti lidského a strojového tlumočení.

V Praze dne 29. 8. 2022

Podpis: