

Prakses iespēja Tildē: Mašīntulkošanas rezultātu kvalitātes cilvēkvērtēšanas tiešo vērtējumu rīkkopas izstrāde

Mašīntulkošanas sistēmu izstrādē iespējams svarīgākais posms ir rezultātu vērtēšana. Praksē tiek izmantotas dažādas automātiskās metrikas, taču kā svarīgākais joprojām tiek uzskatīts cilvēku, it īpaši profesionālu tulkotāju vērtējums. Jau ilgāku laiku Tilde MT iekšējos procesos lieto rīku, kas ļauj salīdzināt divu sistēmu tulkojumu savstarpējo kvalitāti – sistēma A labāka par sistēmu B, B labāka par A, vai A un B vienlīdz labas/sliktas.

MACHINE TRANSLATION EVALUATION

Les patients hospitalisés sont attribués au groupe de personnes atteintes de maladies graves en raison de leurs principales maladies, leur âge avancé et des soins de plus en plus complexes fournis par l'hôpital.

- Hospitalizuoti pacientai priskiriami asmenų, sergančių sunkiomis ligomis dėl jų pagrindinių ligų, vyresnio amžiaus ir vis sudėtingesnės liginės teikiamos priežiūros, grupei.
- Ligoniai priskiriami žmonių, sergančių sunkiomis ligomis dėl jų pagrindinių ligų, vyresnio amžiaus ir vis sudėtingesnės liginės teikiamos priežiūros, grupei.
- Undecided / similar

NEXT

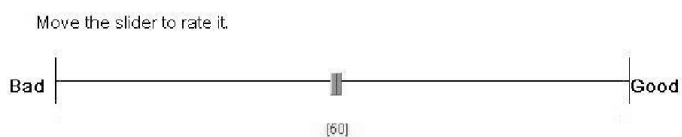
Dear participant of the survey,
We ask you to evaluate our new Machine Translation (MT) system. Usually there are several hundreds of original sentences with two translation variants to be evaluated. You can choose either one of the two translations (strong answer), or choose the "undecided/similar" (weak answer). We expect the minimum of 25 strong evaluation answers from you, and we will appreciate if you do more than that. You can make a break at any time and come back later to continue.

Status of the current evaluation: **incomplete (<25)**

0 sentences evaluated in this survey

Tomēr, šāda rezultātu vērtēšana ir neērta, līdzko jāsāpīdina vairāk nekā divas sistēmas. Turklāt kategoriskie salīdzinājumi "A labāka par B", neļauj atbildēt uz jautājumu, cik daudz viena sistēma ir labāka par otru.

Prakses laikā tiks izstrādāta rīkkopa mašīntulkošanas rezultātu tiešajai vērtēšanai, kurā vērtētājiem m tiks lūgts vērtēt teikuma kvalitāti nepārtrauktā Likerta skalā (https://en.wikipedia.org/wiki/Likert_scale):



Papildus vērtēšanas funkcionalitātei, projektā paredzēts izstrādāt arī rezultātu statistiskās analīzes funkcionalitāti, kas būtu piemērota tiešās vērtēšanas metodoloģijai.

RESULTS OF MACHINE TRANSLATION EVALUATION

Logged in as [REDACTED]

Source - target language: EN - LT

Description:	[REDACTED]
Date created:	9/9/2020 4:01:13 PM
Now:	11/9/2020 5:20:23 PM
Is quantitative:	YES
Escape special chars:	YES

System 1

Name:	[REDACTED] NMT
Url:	[REDACTED]
Bleu score:	23.06

System 2

Name:	Google translation
Url:	https://translate.google.com/
Bleu score:	21.90

Calculation results

Conclusion:	System 1 is better according to the normal approximation interval
Credibility:	Weakly sufficient
Fleiss' kappa interpretation:	0.331 (Fair agreement)
Free kappa interpretation:	0.338 (Fair agreement)
Export data:	CSV

ALL EVALUATIONS

	System 1:	System 2:	Tie:	Total evaluations:
Name	159	119	146	424
p ± ci	37.50% ± 4.61%	28.07% ± 4.28%	34.43% ± 4.52%	
lower	32.89%	23.79%	29.91%	
upper	42.11%	32.34%	38.96%	



All evaluations ignoring ties

Name	p ± ci	lower	upper
System 1	57.19% ± 5.82%	51.38%	63.01%
System 2	42.81% ± 5.82%	36.99%	48.62%



Rīkkopa tiks realizēta kā tīmekļa pakalpojums (web service), kur augstāk minētā funkcionalitāte pamatā tiks īstenota serverpusē. Tiks izstrādāta vienkārša funkcionāla lietotāja saskarne tīmekļa pārlūkprogrammā. Serverpuses kods būs jāizstrādā C# vai Python programmēšanas valodā (atkarībā no kandidāta vēlmēm varam izskatīt arī F#). Kandidātam ir vēlamas zināšanas kādā no šīm valodām.

Kandidātiem būs jāpiedalās darba intervijā prasmju noteikšanai. Prakse notiks sabiedrības Tilde Mašīntulkošanas Grupā.

Praksi vadīs Rihards Krišlauks (rihards.krislauks@tilde.lv), konsultēs Toms Bergmanis un Mārcis Pinnis.