



**NEIRONU MAŠĪNTULKOŠANA MAZAJĀM
VALODĀM – ATSKATS UZ 2018. GADU**

Mārcis Pinnis



21.02.2019.

NEIRONU MAŠĪNTULKOŠANA **Tildē MAZAJĀM**
VALODĀM – ATSKATS UZ 2018. GADU



Mārcis Pinnis

21.02.2019.

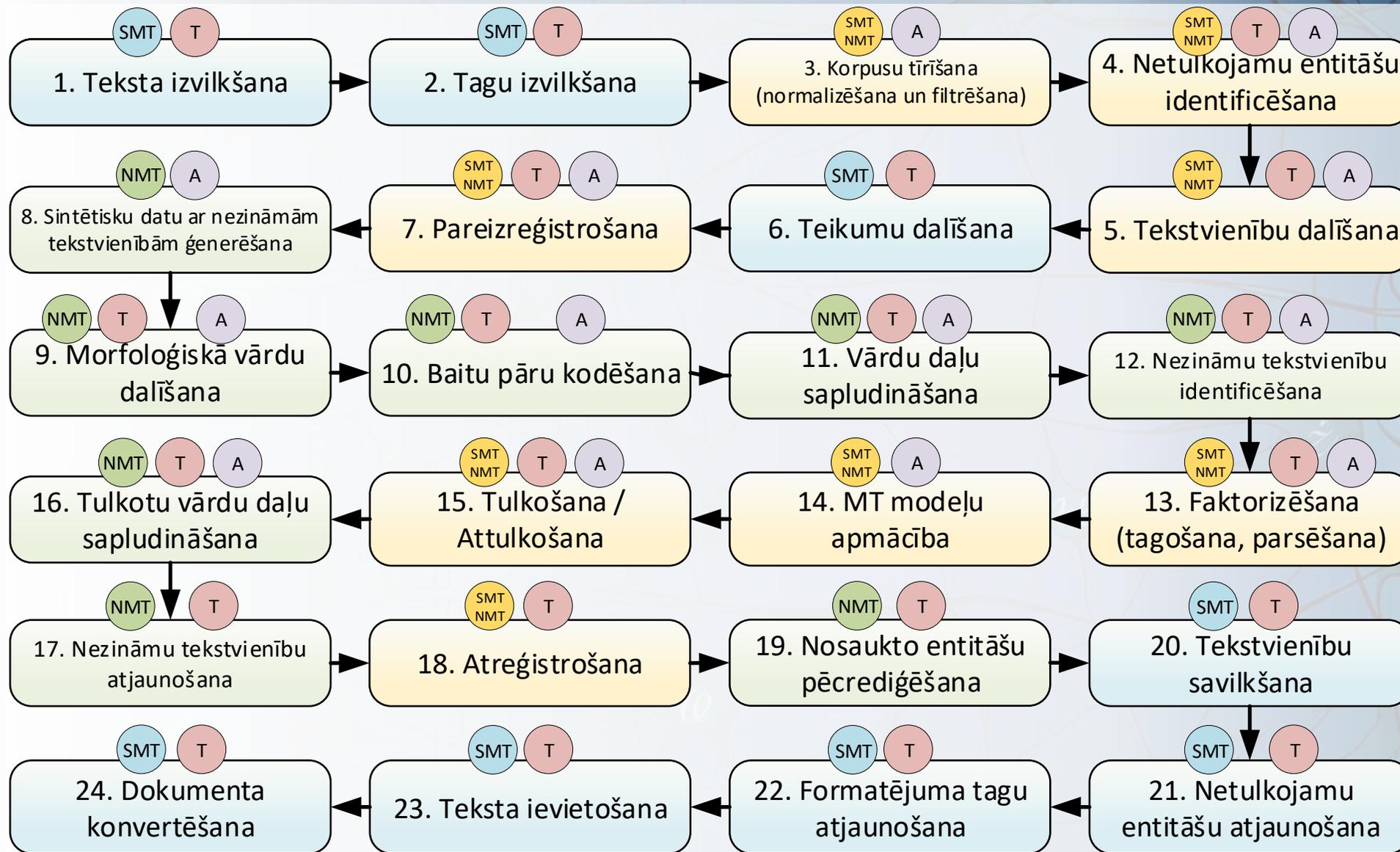
Saturs

- Metodes
- Dalība Mašīntulkošanas konferences sacensībās
- Risināmās problēmas

Metodes

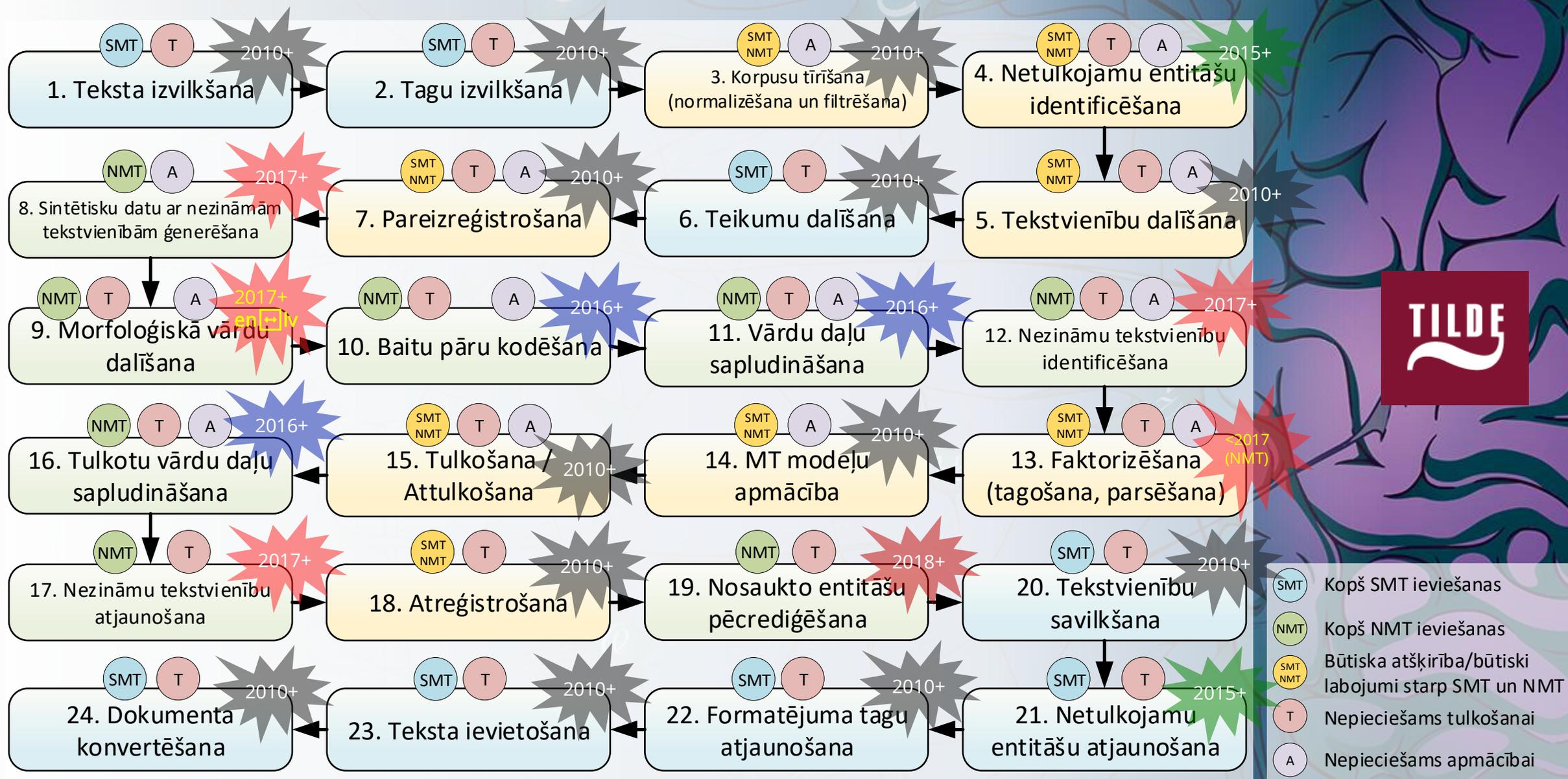


MT darbplūsmas (Tilde MT)

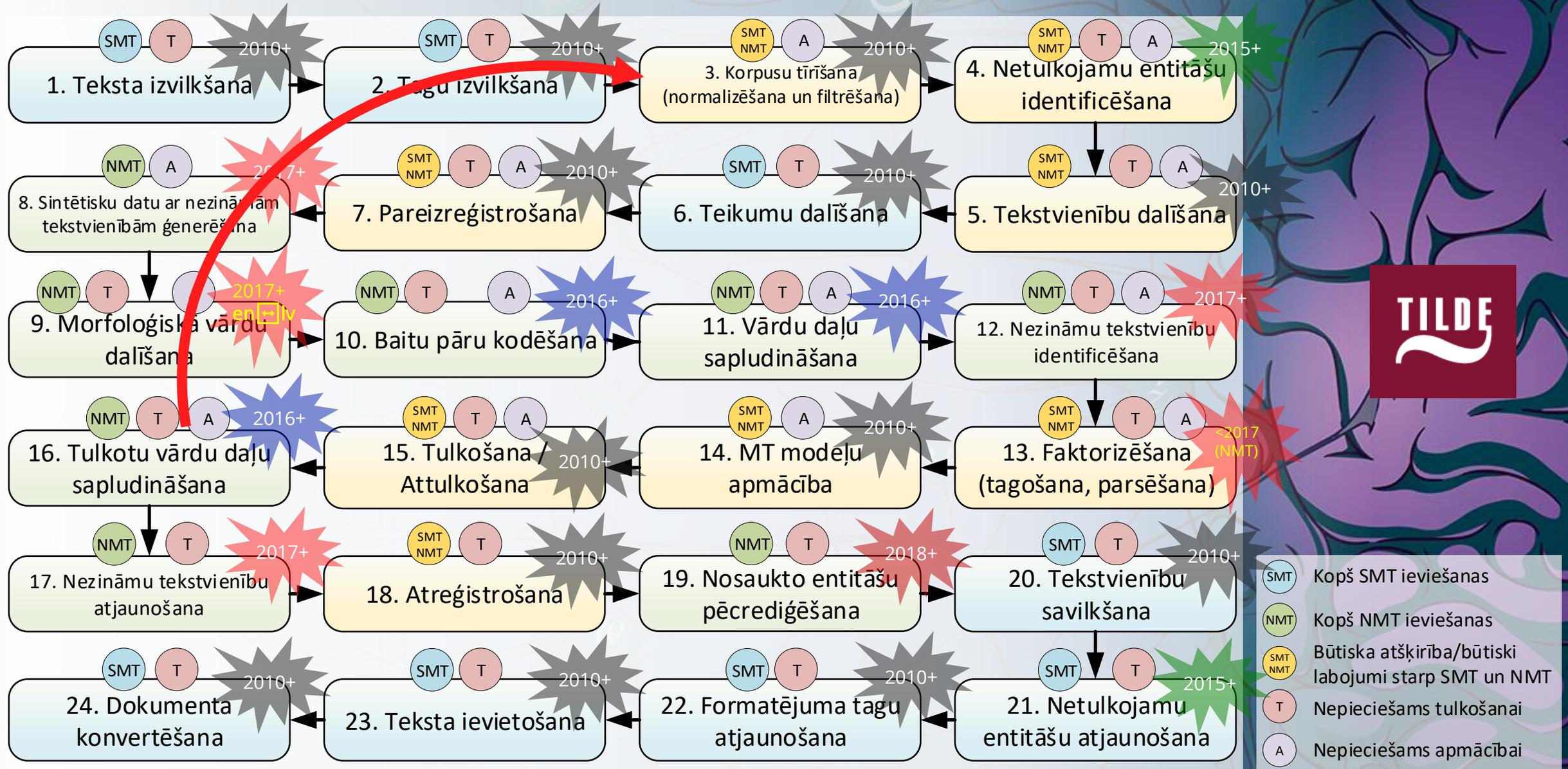


- Kopš SMT ieviešanas
- Kopš NMT ieviešanas
- Būtiska atšķirība/būtiski labojumi starp SMT un NMT
- Nepieciešams tulkošanai
- Nepieciešams apmācībai

MT darbplūsma (Tilde MT)



MT darbplūsuma (Tilde MT) ar attulkošanu



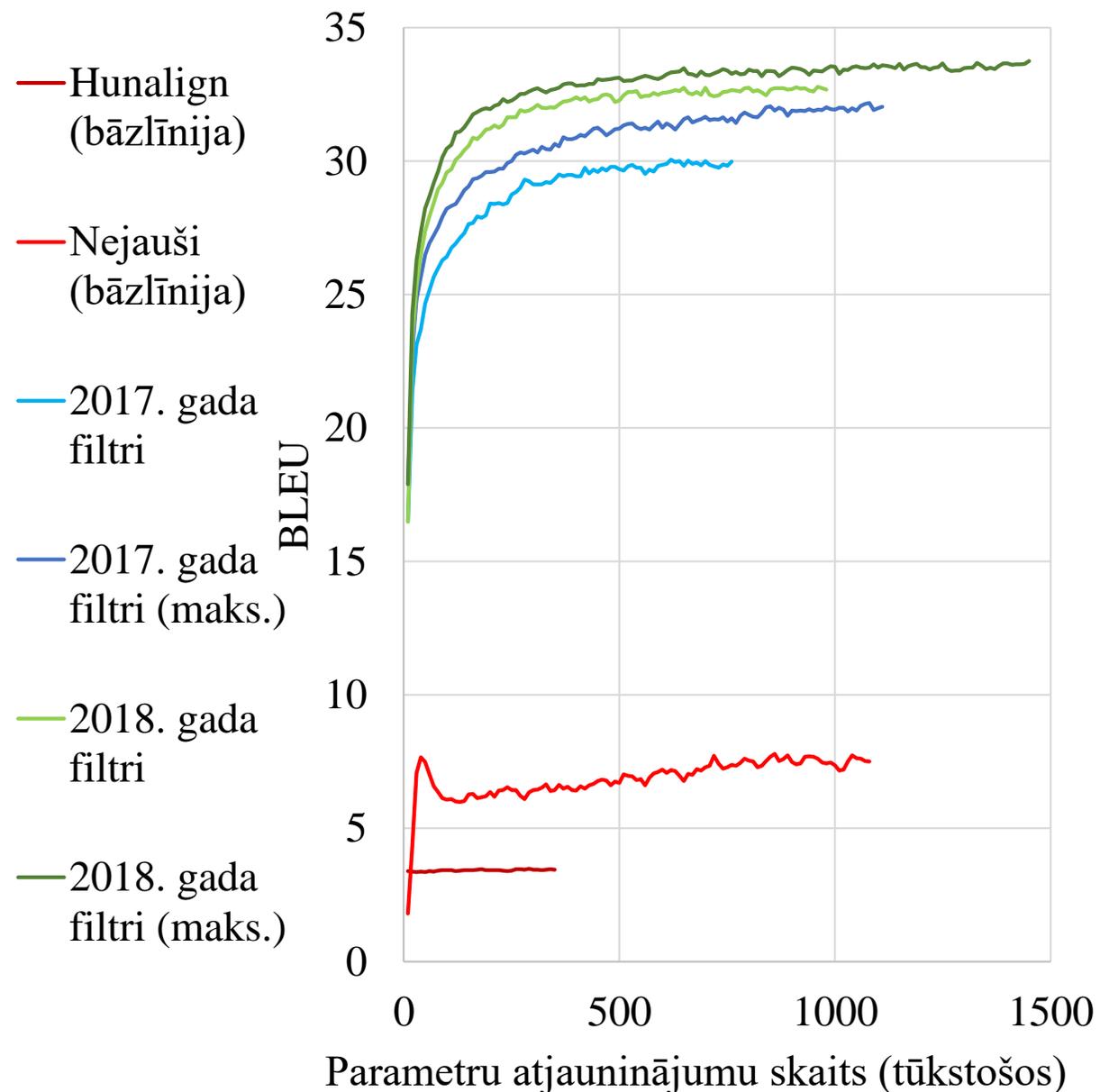
Datu tīrīšana

Datu kvalitāte (tīrība) ir svarīga situācijās, kad ir vēlme iegūt lietojamu (konkurētspējīgu) rezultātu.

Filtru uzdevums - atrast un izdzēst:

- neparalēlus datus;
- datus nepareizās valodās;
- sabojātus datus.

Piemērs no vācu-angļu datu filtrēšanas uzdevuma WMT 2018

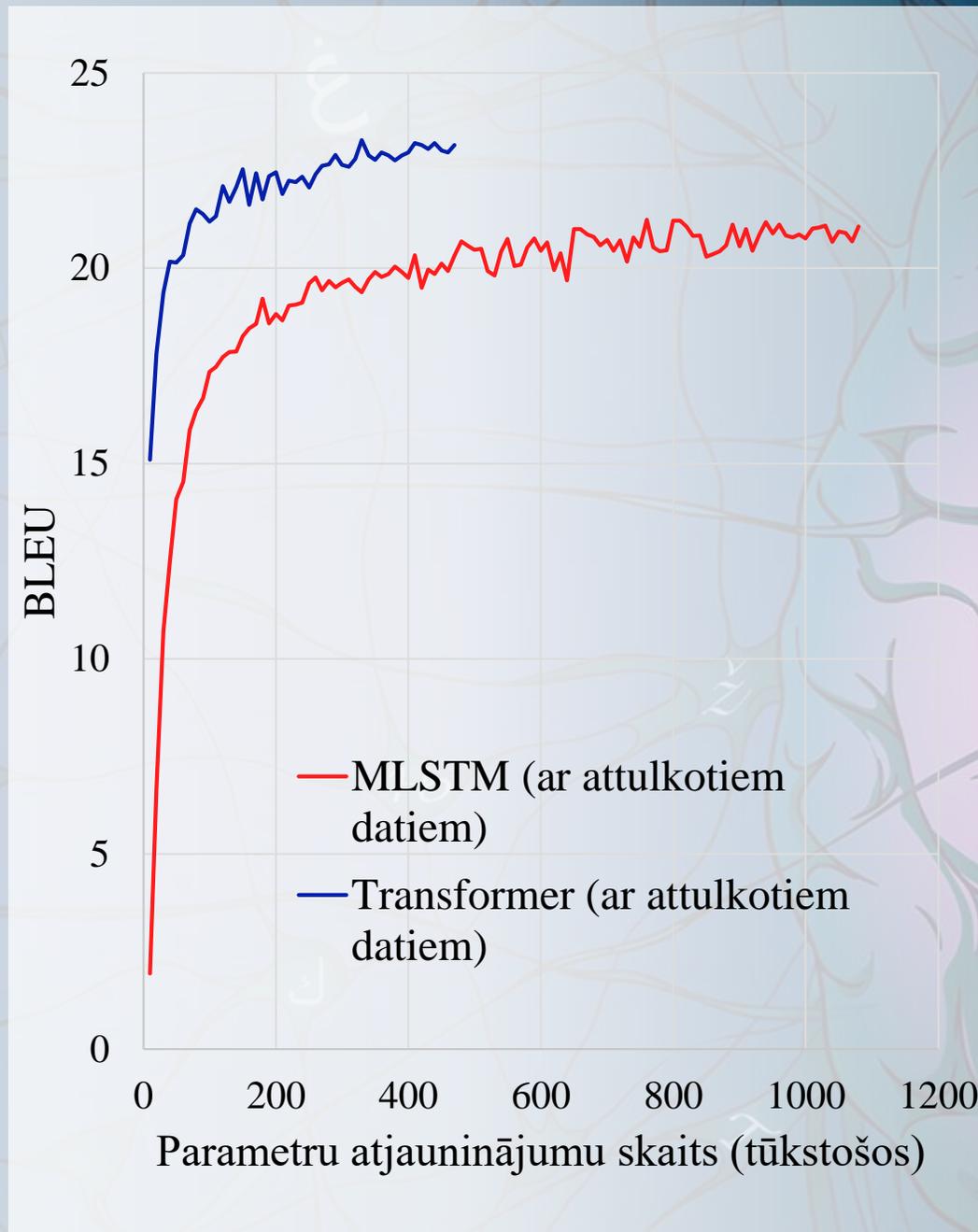


Neironu tīklu arhitektūras

2017. gadā - MLSTM
(*multiplikatīvās garās
īstermiņa atmiņas
rekurentie neironu tīkli*)

2018. gadā - Transformer
(*jeb «pašuzmanības»
neironu tīkli.*)

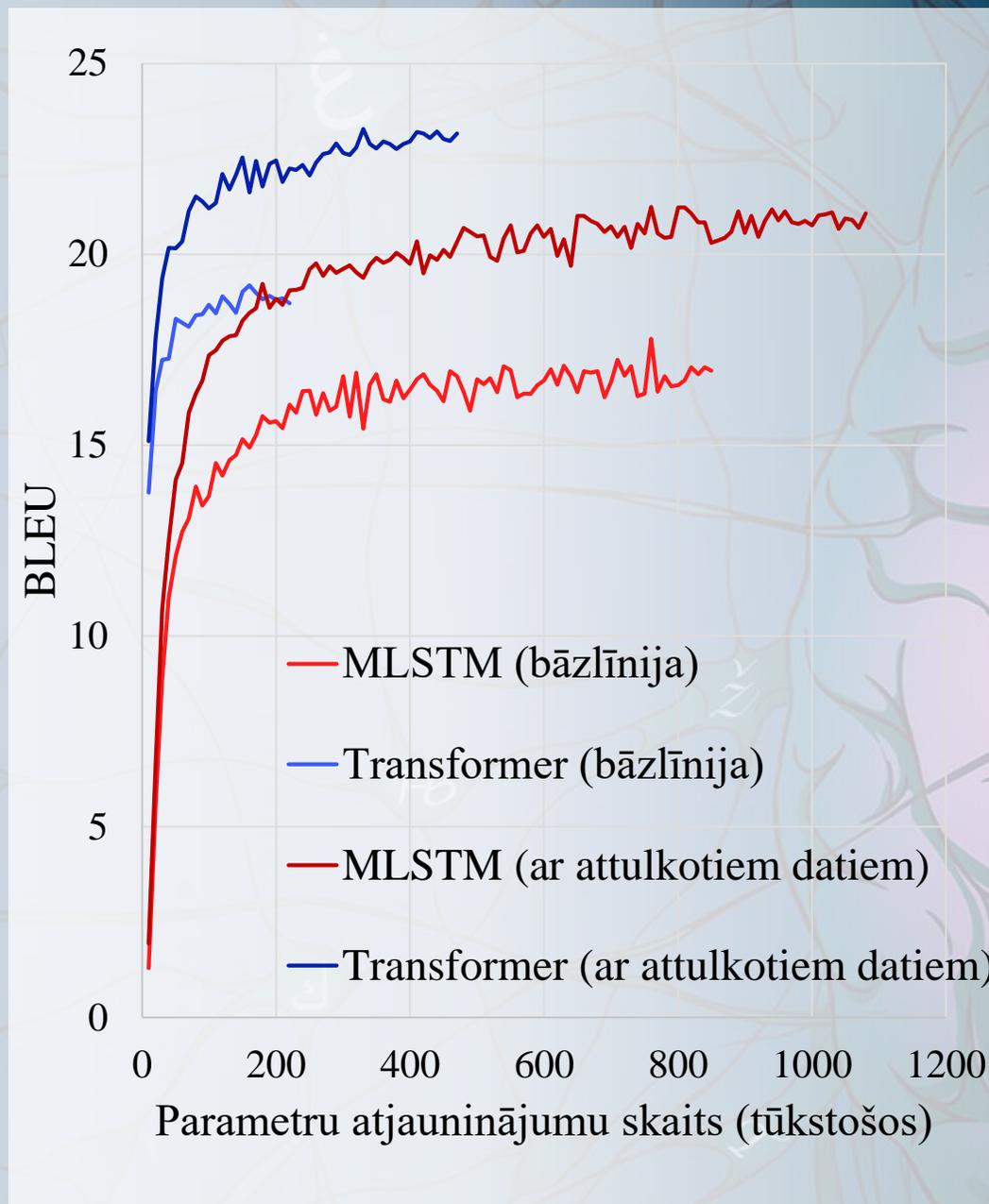
*Piemērs no Tildes angļu-
igauņu WMT 2018
eksperimentiem*



NMT pielāgošana jomai

Attulkošana - efektīvākā
NMT sistēmu
pielāgošanas metode,
situācijās, kad jomas
paralēlo datu nav.

*Piemērs no Tildes angļu-
igauņu WMT 2018
eksperimentiem*



Kur valodu specifika?

2017. gadā mēs pētījām:

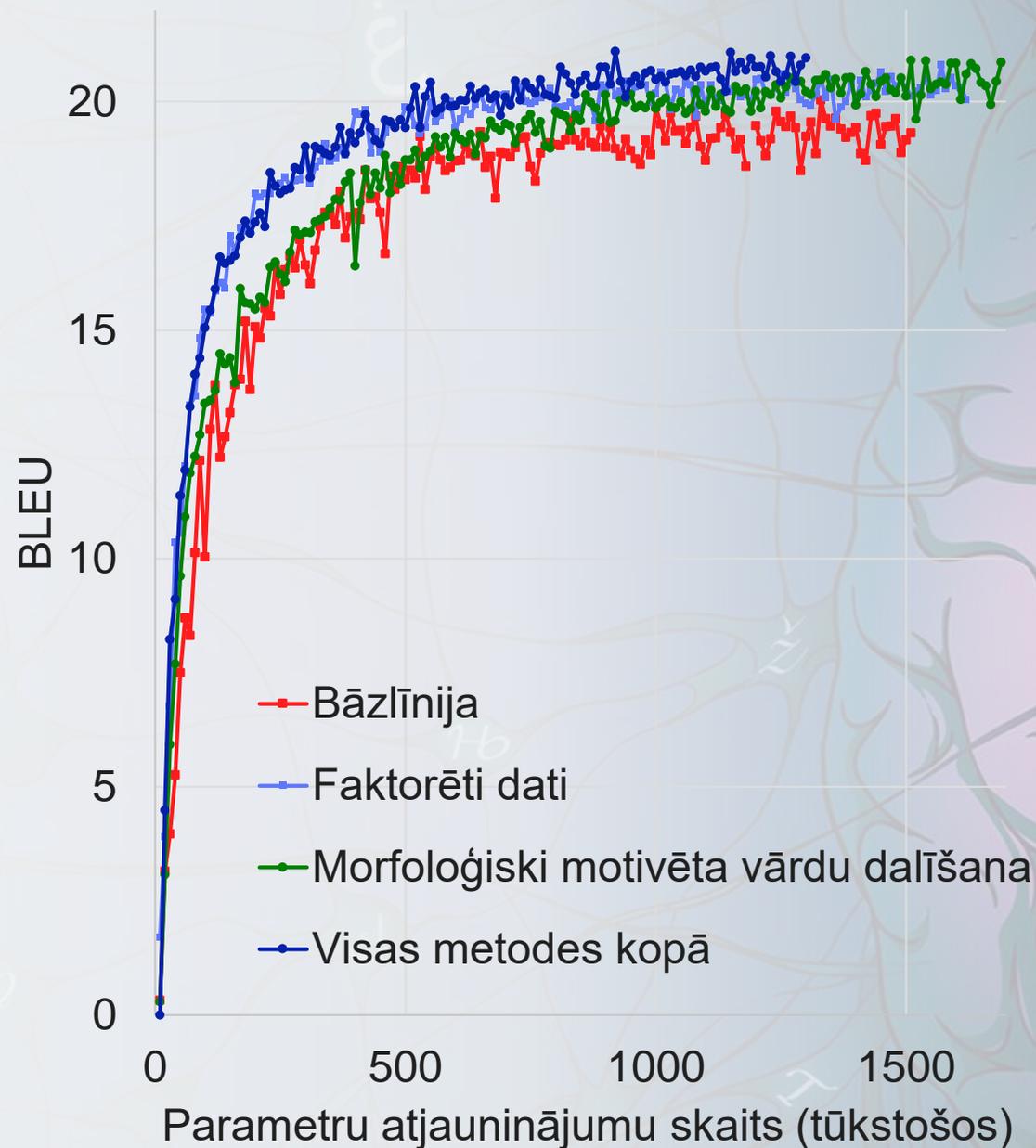
- morfoloģiski motivētu vārdu dalīšanu
- faktorētu datu izmantošanu

Ieguvumi - ~+1 BLEU

Tīrāku datu un attulkošanas (labākas pielāgošanas) ieguvums ir lielāks (≥ 5 BLEU)

Jaunāku modeļu ieguvums lielāks (≥ 2 BLEU)

Vispirms jāatrod labākās no valodām neatkarīgās metodes...



Dalība Mašīntulkošanas konferences sacensībās



Sacensības

Uzdevums: izmantojot organizatoru sagatavotus datus (*ierobežotu datu scenārijs*) vai citus datus (*neierobežotu datu scenārijs*), izstrādāt mašintulkošanas sistēmu ziņu tulkošanai, ar kuru ir jāpārtulko novērtēšanas datu kopa; pārtulkotā datu kopa ir jāiesniedz vērtēšanai.

Dalībnieki: mašintulkošanas jomas pētnieki un izstrādātāji no pētniecības centriem un uzņēmumiem no visas pasaules

Sacensības

Vērtēšana: Tulkojumi tiek vērtēti, izmantojot:

- 1) *automātiskas vērtēšanas metodes,*
- 2) *manuālas vērtēšanas metodes, iesaistot cilvēkus* (tulkotājus, redaktorus, valodu ekspertus, kā arī «pūli»).

Cilvēkvērtējums tiek uzskatīts par uzticamāku vērtējumu!



The screenshot shows the 'Appraise' interface for WMT 2018. At the top, it displays 'Appraise Dashboard' and the user 'engest5011'. Below this, a progress bar indicates '0/10 blocks, 10 items left in block' and the task ID 'WMT18RefDA #611:Segment #1891'. The language pair is 'English → Estonian (eesti)'. The source text is: 'Ilmselt ikka selleks, et rikkamad kliendid ei peaks end lihtsurelikena tundma, kes mingite odavate kottidega ringi promeneerivad.' The candidate translation is: 'Pole kahtlust, et jõukamad kliendid ei pea tundma, et nad kogu telefoni esipaneeli, kaotatud on odava kotiga ringi liiguvad.' A slider below asks 'How accurately does the above candidate text convey the original semantics of the source text? Slider ranges from Not at all (left) to Perfectly (right)'. The slider is currently positioned at approximately 25%. There are 'Reset' and 'Submit' buttons at the bottom.

Attēls no WMT 2018 mašīntulkošanas sistēmu vērtēšanas sistēmas (<http://wmt18.appraise.cf>)

س

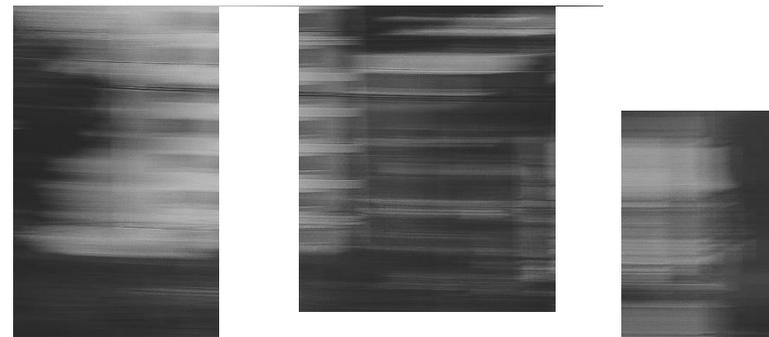
Sacensības

Mēs sacensībās piedalījāties:

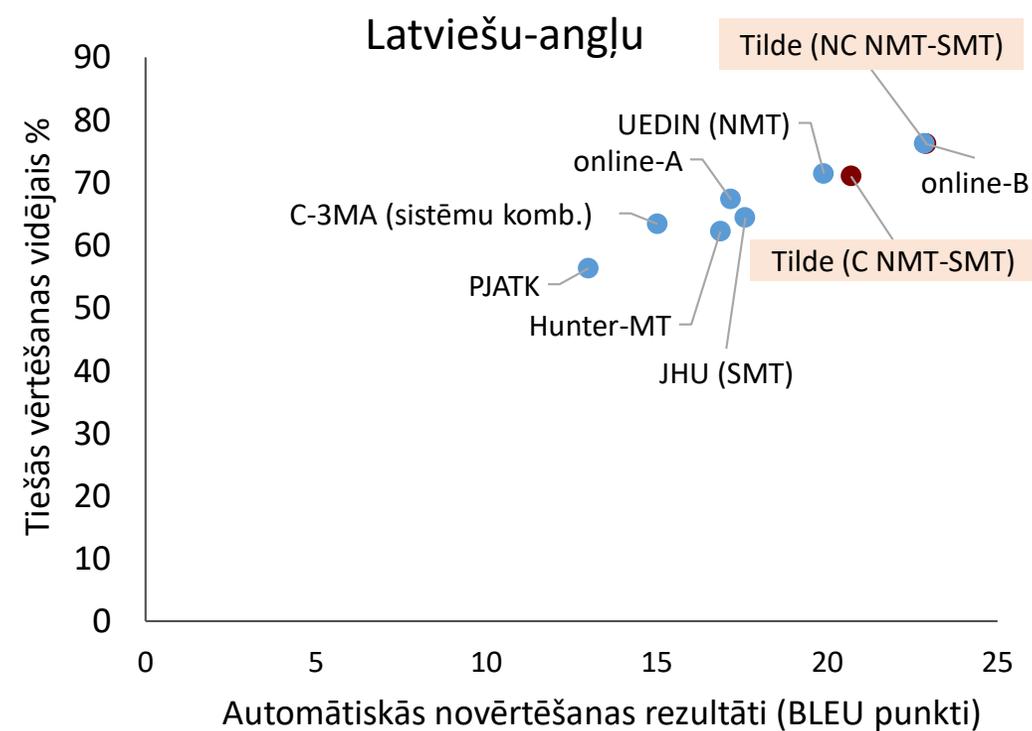
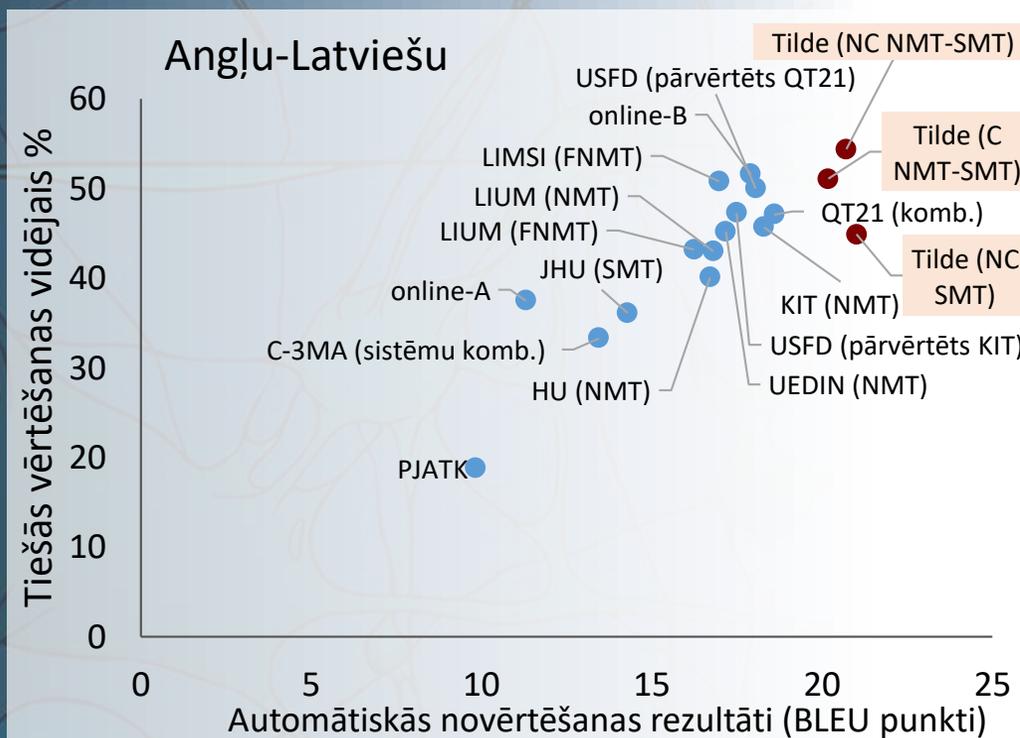
- 2017. gadā, izstrādājot **angļu-latviešu** mašīntulkošanas sistēmas ziņu jomai
- 2018. gadā, izstrādājot **angļu-igauņu** mašīntulkošanas sistēmas ziņu jomai

L

ق

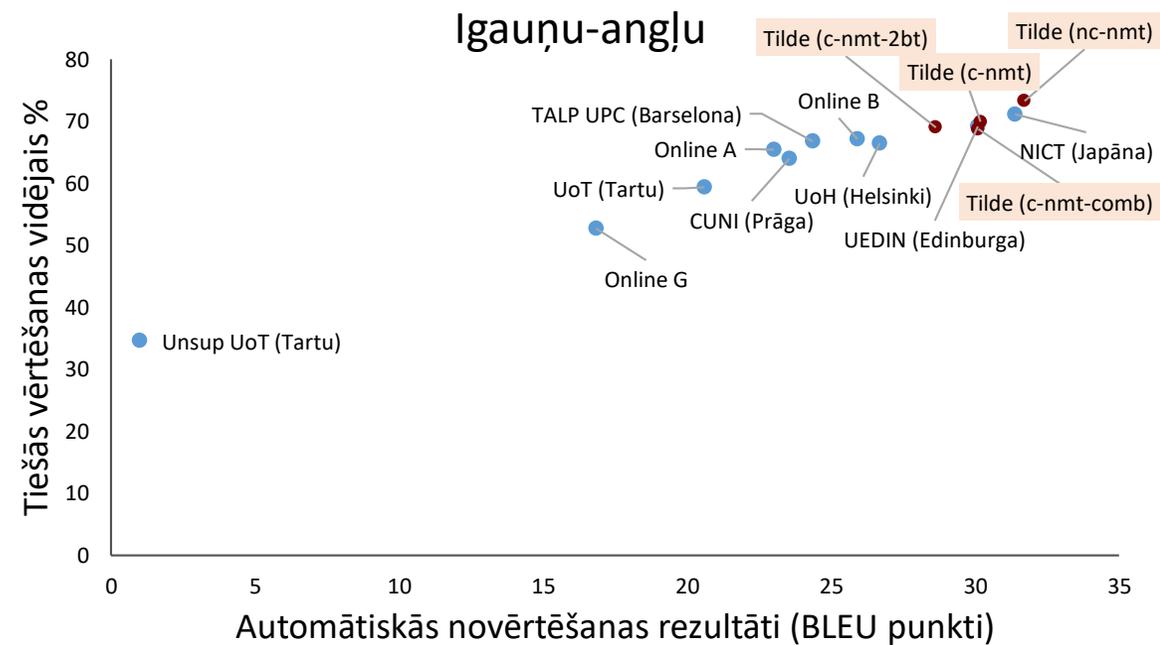
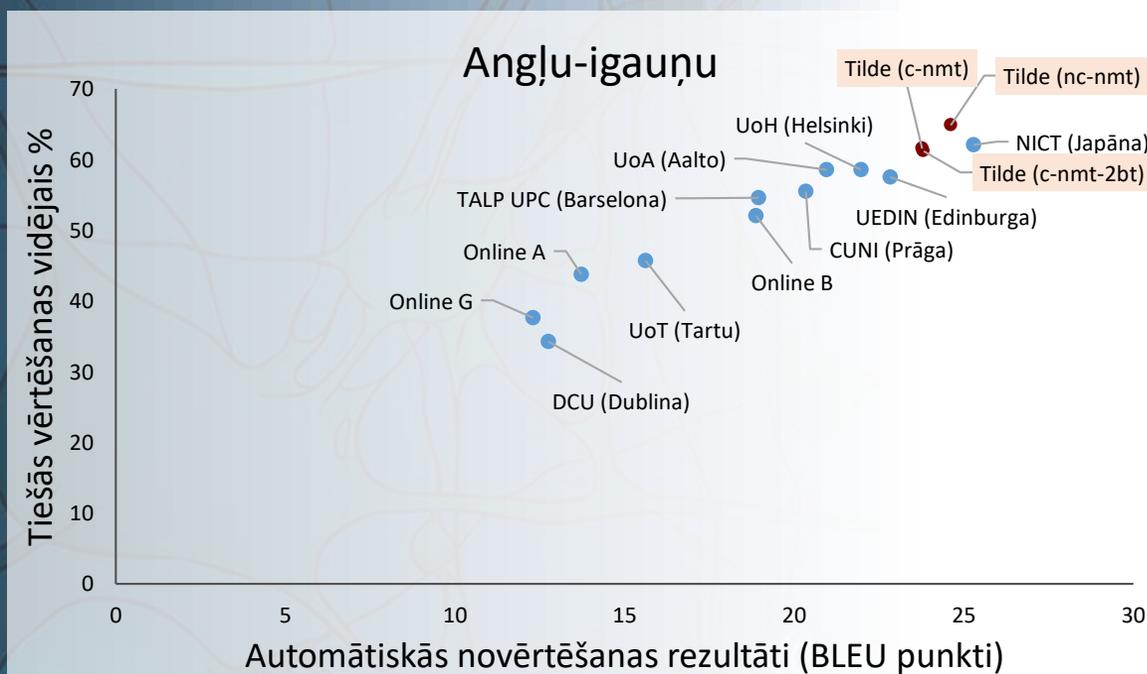


Vislabākās angļu↔latviešu neironu mašintulkošanas sistēmas WMT 2017 sacensībās



Mašintulkošanas sistēmas apmācītas ar tā laika spējīgākajiem neironu tīklu modeļiem mašintulkošanā - uzmanībās balsītiem rekurentiem neironu tīkliem ar multiplikatīvās garās īstermiņa atmiņas šūnām rekurentajos slāņos.

Vislabākās angļu↔igauņu neironu mašintulkošanas sistēmas WMT 2018 sacensībās



Mašintulkošanas sistēmas apmācītas, izmantojot 2018. gada spējīgākos modeļus - Transformer modeļus.

(Dažas) problēmas

**...jeb aktuālas tēmas šogad
un tuvākā nākotnē**



Retu vārdu/frāžu pareizāka tulkošana: **nosauktās entitātes**

Latviešu	Angļu
Google NMT (2019. gada 20. februāris)	
Šuplinska solīja, ka algas [...]	Šuplinska promised to pay [...]
[...] priekšsēdētāja Inga Vanaga.	[...] President Inga Vanaga.
Rīgas domes priekšsēdētājs Nils Ušakovs (S) [...]	Riga City Council Chairman Nils Ushakov (S) [...]
Ansis Ansbergs (LA) interesējās, [...] izpildītājs Emīls Jakrins atbildēja [...]	Ansis Ansberg (LA) was interested in [...] Emperor Yakrins answered [...]
Tilde NMT (WMT 2018)	
Šuplinska solīja, ka algas [...]	Shuplinska promised that wages [...]
[...] priekšsēdētāja Inga Vanaga.	[...] Chairman Inga Vanaga.
Rīgas domes priekšsēdētājs Nils Ušakovs (S) [...]	Riga City Council Chairman Nils Ushakov (S) [...]
Ansis Ansbergs (LA) interesējās, [...] izpildītājs Emīls Jakrins atbildēja [...]	Ansis Ansbergs (LA) was interested, [...] performer Emile Yakrin replied [...]



Retu vārdu/frāžu pareizāka tulkošana: **terminoloģija**

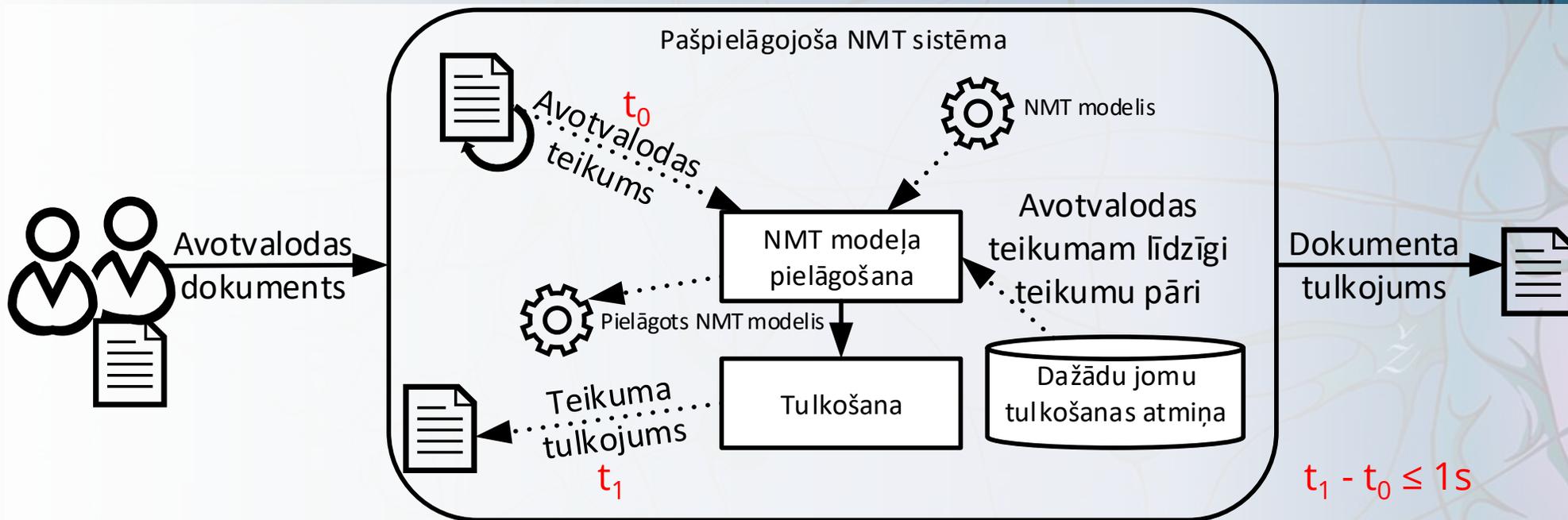
Latviešu	Angļu
<i>Google NMT (2019. gada 20. februāris)</i>	
The corncrake , corn crake or landrail is a bird in the rail family .	Kukurūzas , kukurūzas krējums vai krastmala ir putns dzelzceļa ģimenē .
<i>Tilde NMT (WMT 2018)</i>	
The corncrake , corn crake or landrail is a bird in the rail family .	Kukurūzas krakšķis , kukurūzas krakšķis vai zemūdene ir putnu dzimtas putns .

corncrake, corn crake, landrail
bird
rail family

grieze
putns
dumbrvistiņu dzimta



Pielāgošana jomai - pašpielāgojošas NMT sistēmas



Sākotnējie rezultāti liecina, ka plašas jomas sistēmām, ieviešot pašpielāgojošas NMT sistēmas iespējas, tulkojumu kvalitāti var palielināt **no +2** (ja līdzīgu datu tulkošanas atmiņā nav daudz) **līdz pat +50** (ja līdzīgi jomas dati ir tulkošanas atmiņā sastopami) **BLEU**

Kā gūt labumu no (daudz) datiem?

Tulkošanas virziens	Datu kopapjoms (teikumu pāri milj.)	BLEU	Tiešā vērtējuma vidējais svērtais rezultāts
Angļu-igauņu (WMT 2018)	3,65	23,54	61,6
	63,07	24,35	64,9
Igauņu-angļu (WMT 2018)	3,73	29,46	69,9
	67,91	30,94	73,3
Angļu-latviešu (WMT 2017)	6,19	20,18	51,1
	46,04	20,71	54,4
Latviešu-angļu (WMT 2017)	6,19	20,81	71,0
	45,71	23,02	76,2

Pašpielāgojošas NMT sistēmas?



Paldies!

