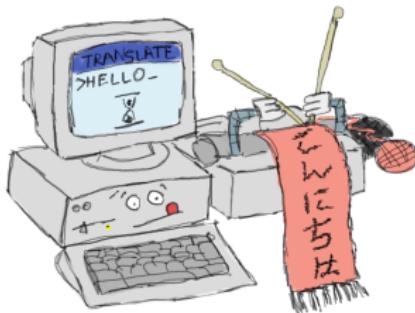


# Do hlubin překladače Charles Translator

Mgr. Martin Popel, Ph.D.

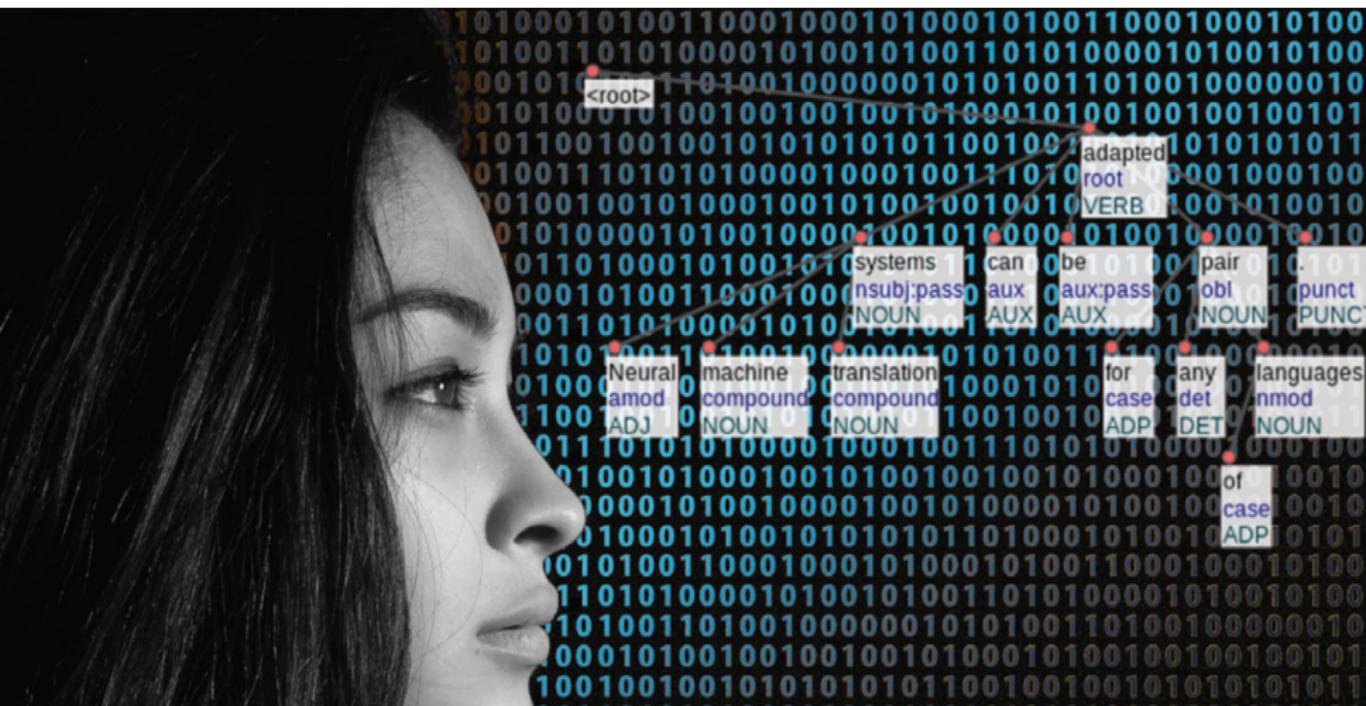
ÚFAL (Institute of Formal and Applied Linguistics), MFF UK

2024-03-13, Matematické problémy nematematiků



source	Great talkers are little doers.
Yandex	Velké talkers jsou trochu činitelé.
Bing	Velcí vysílačky jsou malí činitelé.
Google	Velcí mluvčí jsou malí lidé.
TectoMT	Velcí řečníci jsou malí vrazi.
CUBBITT	Velcí mluvkové jsou malí dřiči.

# Natural Language Processing?



## Spell check vs. grammar check

3

Chlapci šly.  
Chlapec šli do školy.

## Spell check vs. grammar check

3

Dívce nešly hodinky. Chlapci šly.  
Chlapec šli do školy.

# Spell check vs. grammar check

3

Dívce nešly hodinky. Chlapci šly.  
Kdo kam co donesl? Chlapec šli do školy.



# Spell check vs. grammar check

3

<http://ufal.cz/korektor>

Korektor-sample.rtf — Edited

Potkávám je na každém **krou**.

Did you mean kroku?

Lední medvěd byl schován za **krou**.

Vláďa nás **přivíral** s otevřenou **náruči**.

Did you mean přivítal? Did you mean náručí?

Kocour **přivíral** oči slastí, když pozoroval kanárka v **klexi**.

Tento víkend raději zůstaneme v **praze**.

Error in capitalization? Did you mean Praze?

Byl by to rytíř, kde v pláně hřích vzlet,  
Vědě jsem jse seheldo na přídoutně v světě si nezastavá:  
„Ukryjemné, chvěla, milý nás jest

Kolem jsou jest vyhrávaných  
A svítí co pláčem, rád pravil:  
Ale plná jízdo zaporodilo se, vys.

již dávno vás poháru a vlanných rány,  
v jablonění je píše je i v kristování,

srdce v své ženských svém  
v obly pětky tam a vzíti,  
na kóňku je, milý svěžek.

I'll come a bit later on my own.

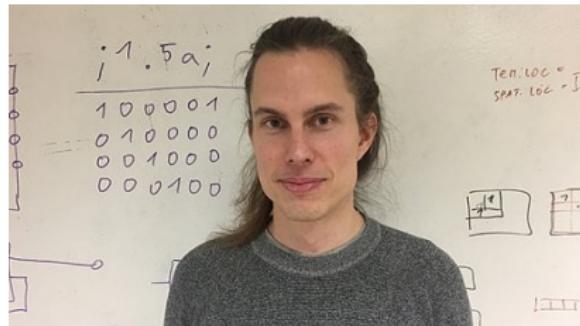
I'll come a bit later on my own.  
Sem čelist ještě na své milé.

Can you add and subtract numbers?  
What about words and images?

king - man + woman = ?

$$\text{king} - \text{man} + \text{woman} = \text{queen}$$

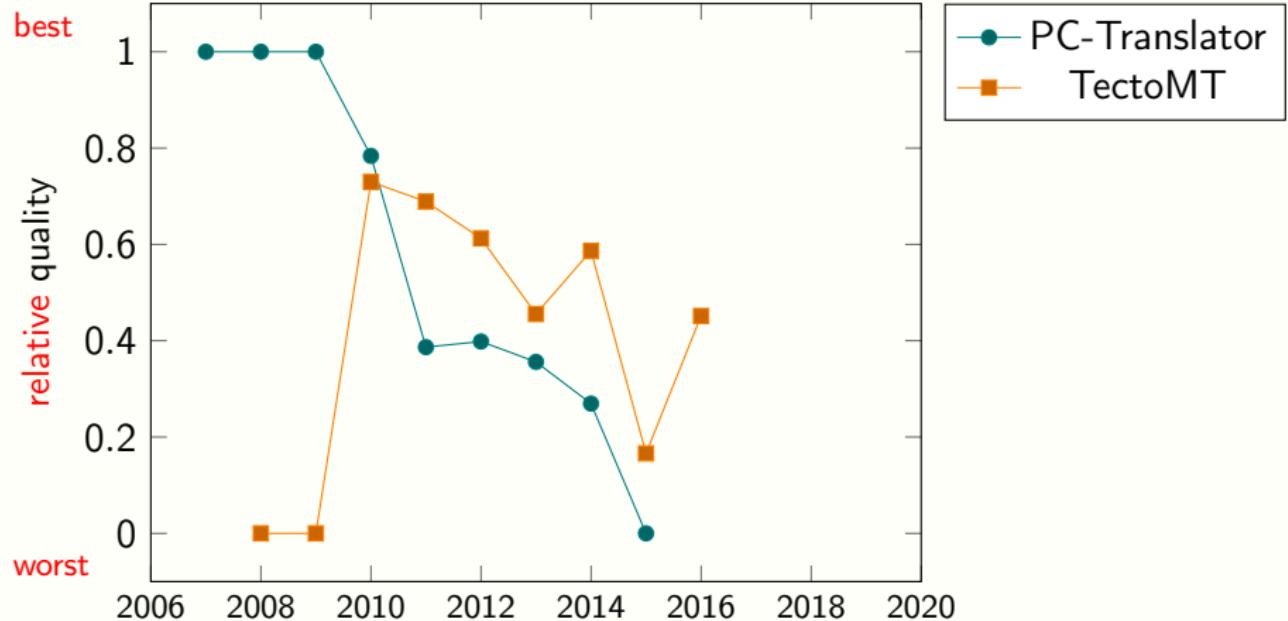
Tomáš Mikolov, 2012, word2vec



<https://projector.tensorflow.org/>

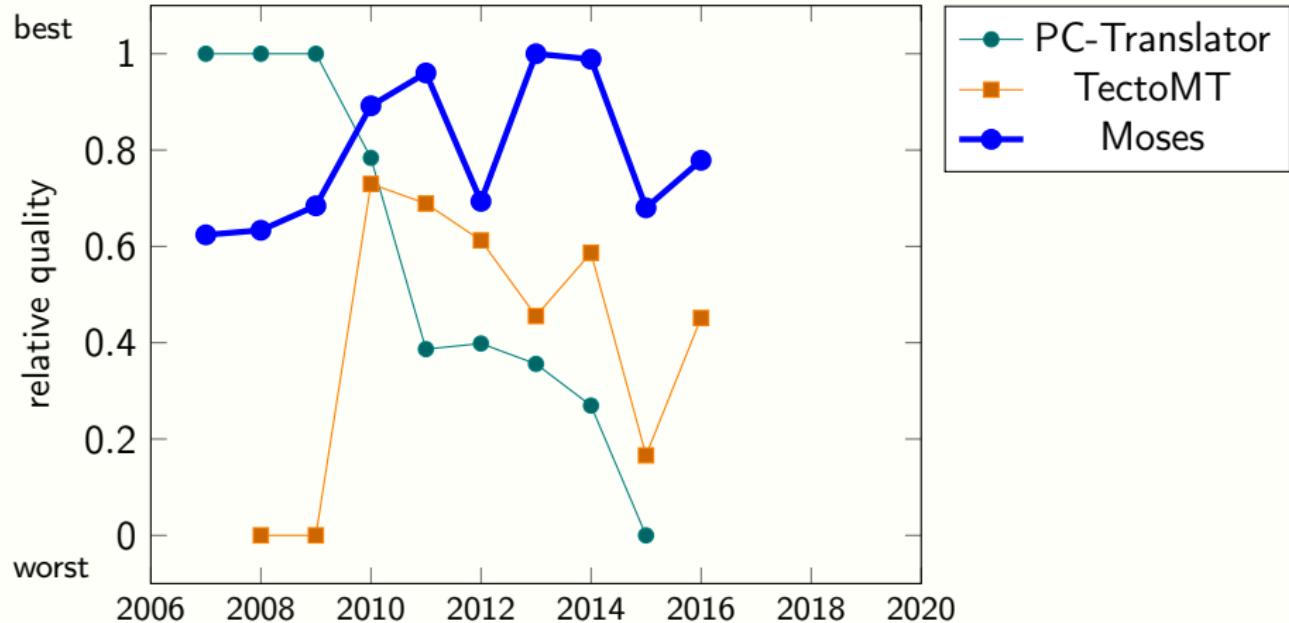
# English→Czech MT in 2007–2020

6



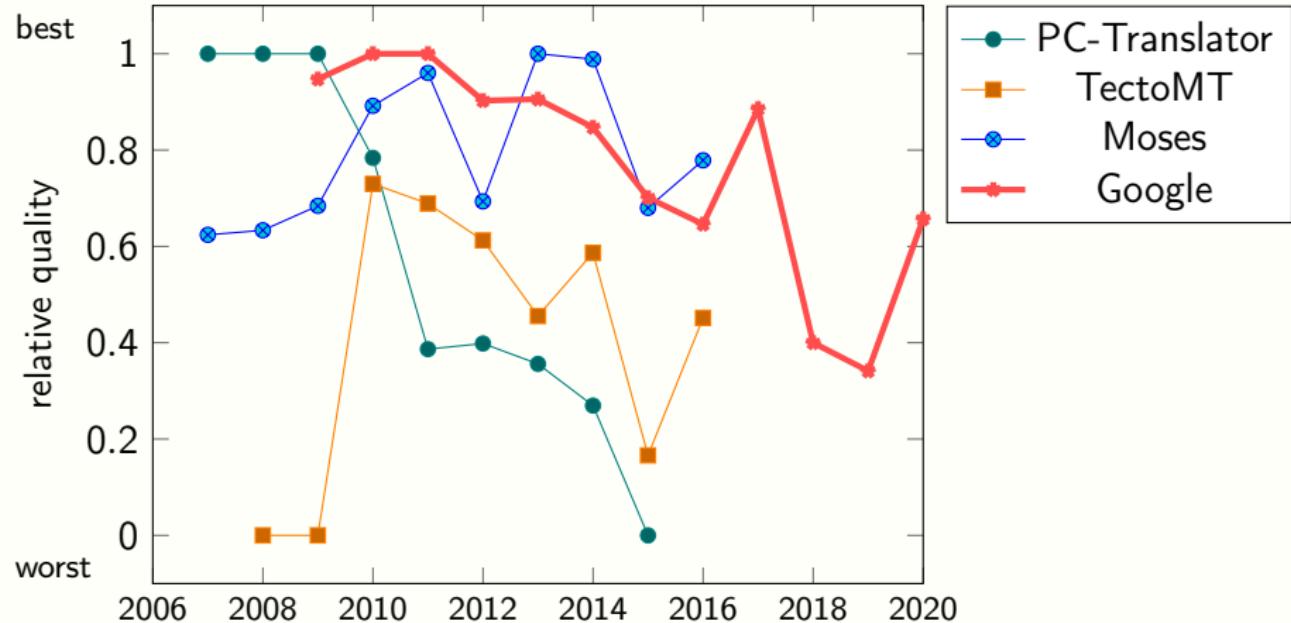
# English→Czech MT in 2007–2020

6



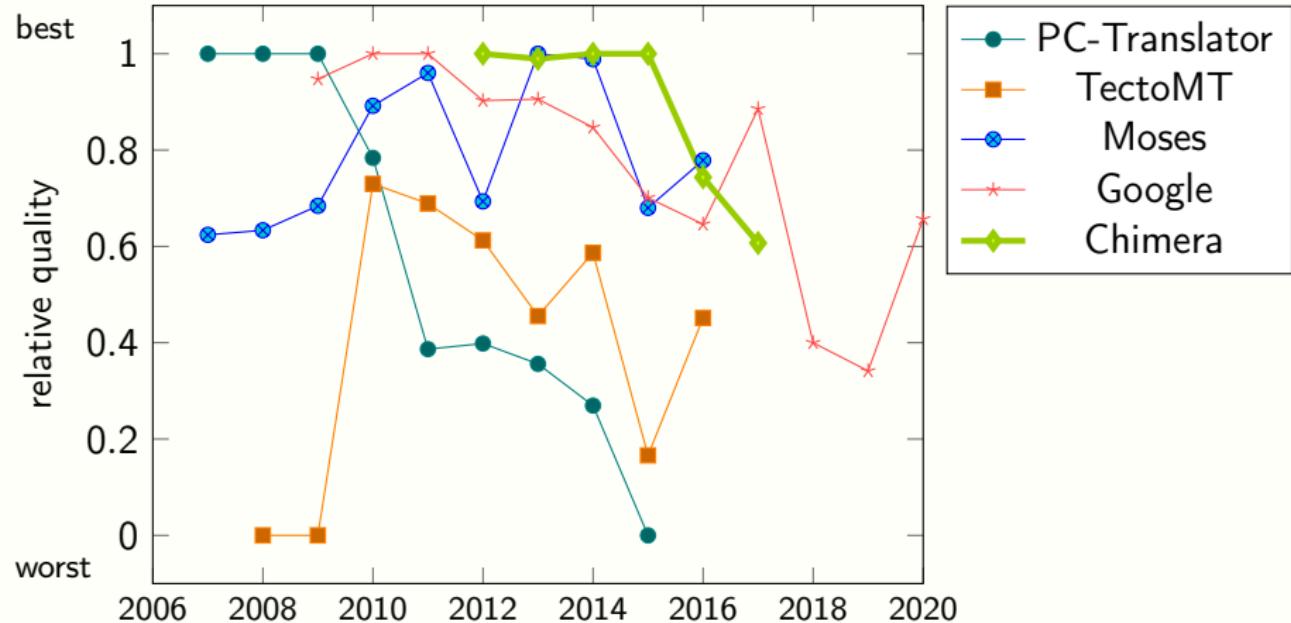
# English→Czech MT in 2007–2020

6



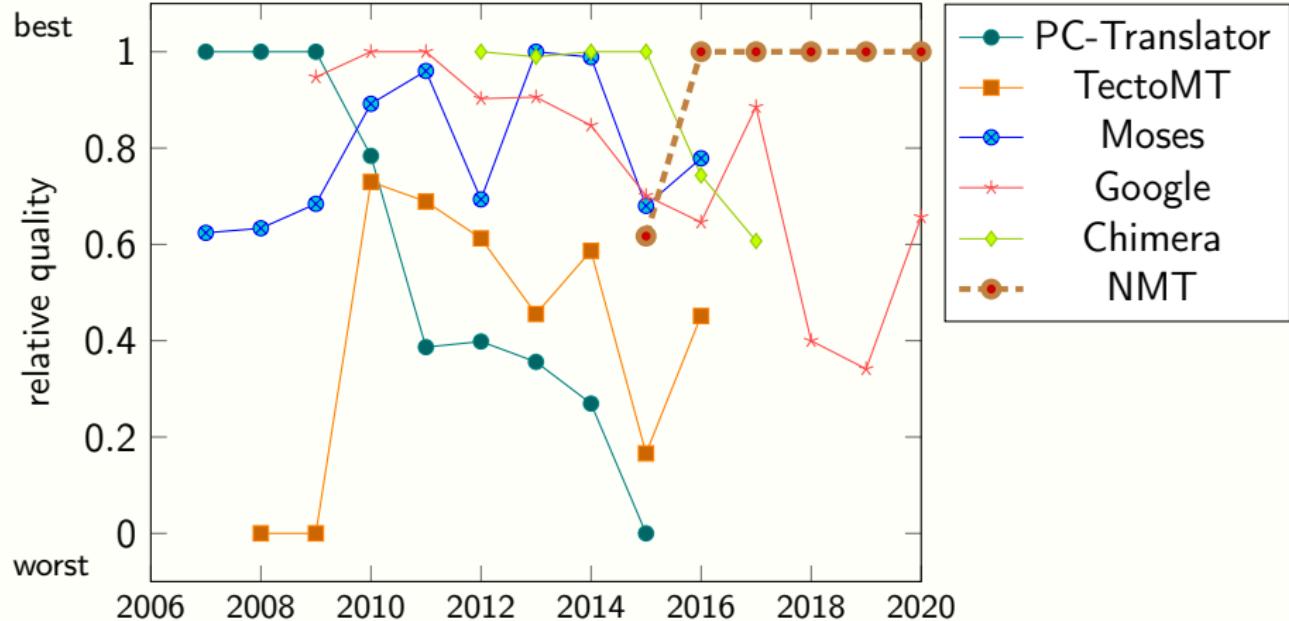
# English→Czech MT in 2007–2020

6



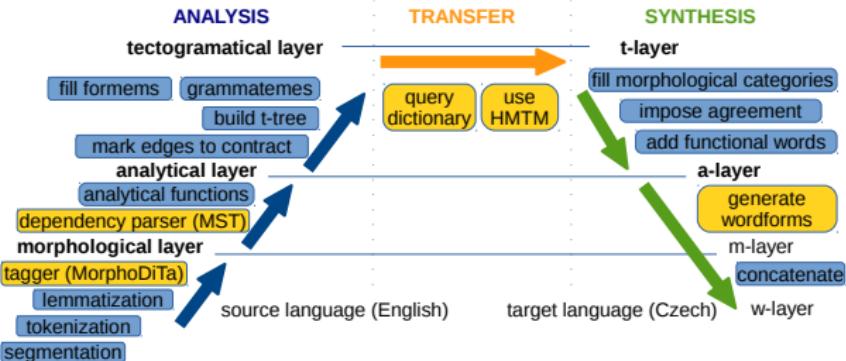
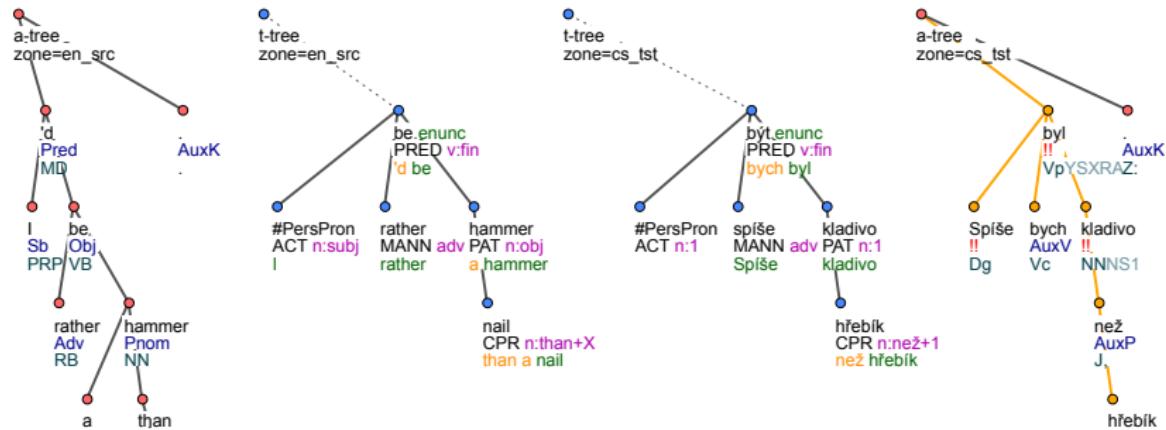
# English→Czech MT in 2007–2020

6



# Deep-syntactic translator TectoMT

7



I'd rather be a hammer than a nail.

Spíše bych byl kladivo než hřebík/nehet.

# Machine learning: features in the translation model

output_label=hřebík#N	
feature	$\lambda$
child_formeme_n:in+X=1	1.64
is_member_of_coord=1	1.30
child_formeme_v:fin=1	1.04
next_lemma=down	0.84
is_capitalized=1	0.79
<b>+precedes_parent=0</b>	<b>0.75</b>
tense_g=post	0.74
<b>+voice_g=active</b>	<b>0.66</b>
prev_lemma=drive	0.66
parent_capitalized=1	0.62
formeme=n:from+X	0.60
<b>+prev_lemma=hammer</b>	<b>0.59</b>
child_lemma_few=1	0.55
child_lemma_remove=1	0.54
sempos=n.denot	0.50
next_lemma=and	0.50
formeme_g=v:until+fin	0.49
child_lemma_rusty=1	0.47
...	

# Machine learning: features in the translation model

output_label=hřebík#N	
feature	$\lambda$
child_formeme_n:in+X=1	1.64
is_member_of_coord=1	1.30
child_formeme_v:fin=1	1.04
next_lemma=down	0.84
is_capitalized=1	0.79
+precedes_parent=0	<b>0.75</b>
tense_g=post	0.74
+voice_g=active	<b>0.66</b>
prev_lemma=drive	0.66
parent_capitalized=1	0.62
formeme=n:from+X	0.60
+prev_lemma=hammer	<b>0.59</b>
child_lemma_few=1	0.55
child_lemma_remove=1	0.54
sempos=n.denot	0.50
next_lemma=and	0.50
formeme_g=v:until+fin	0.49
child_lemma_rusty=1	0.47
...	

output_label=nehet#N	
feature	$\lambda$
child_formeme_n:poss=1	1.32
child_lemma_finger=1	1.07
child_formeme_n:of+X=1	0.98
precedes_parent=1	0.88
prev_lemma=black	0.77
child_lemma_broken=1	0.76
child_formeme_v:attr=1	0.70
formeme=n:at+X	0.67
formeme_g=n:attr	0.67
child_lemma_long=1	0.67
next_lemma=file	0.60
child_lemma_false=1	0.58
prev_lemma=false	0.58
+number=sg	<b>0.56</b>
formeme=n:obj	0.53
formeme=n:by+X	0.52
...	

What is this?

9



40 GPU (GeForce GTX 1080 Ti, 12 billion transistors)

9



# Artificial intelligence and its subfields

10

artificial intelligence

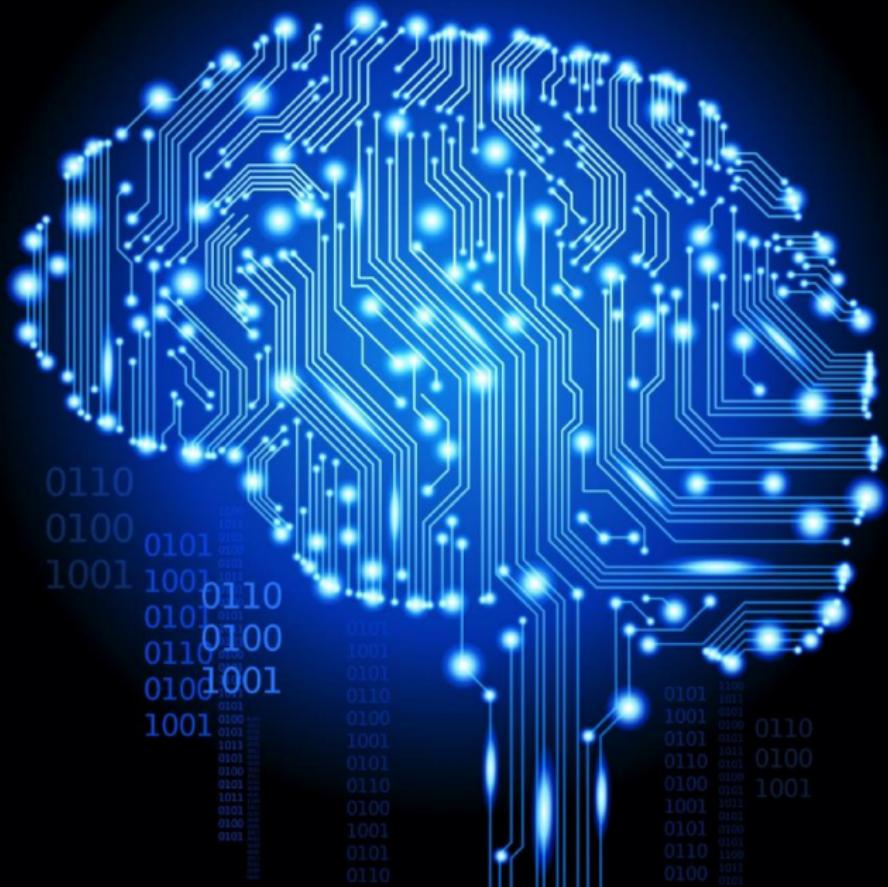
~1950

machine learning

~1980

deep learning

~2010



# Artificial intelligence and its subfields

10

## artificial intelligence

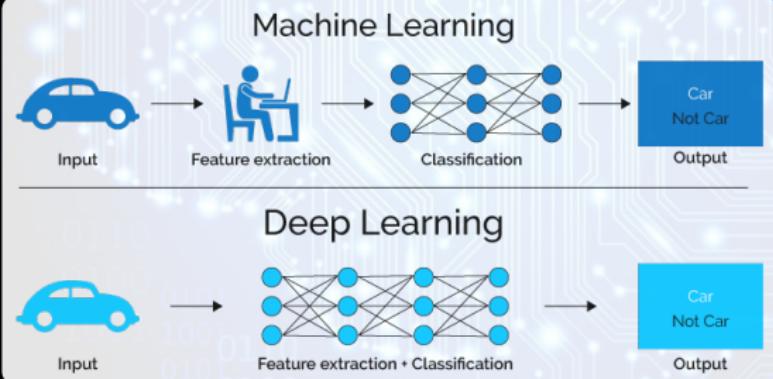
~1950

## machine learning

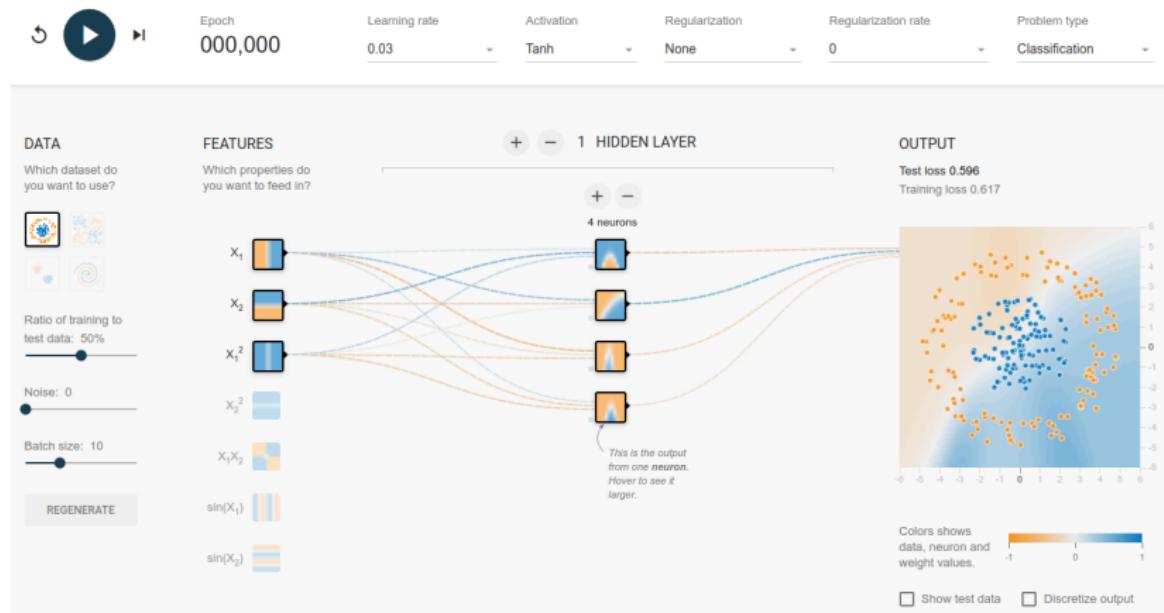
~1980

## deep learning

~2010

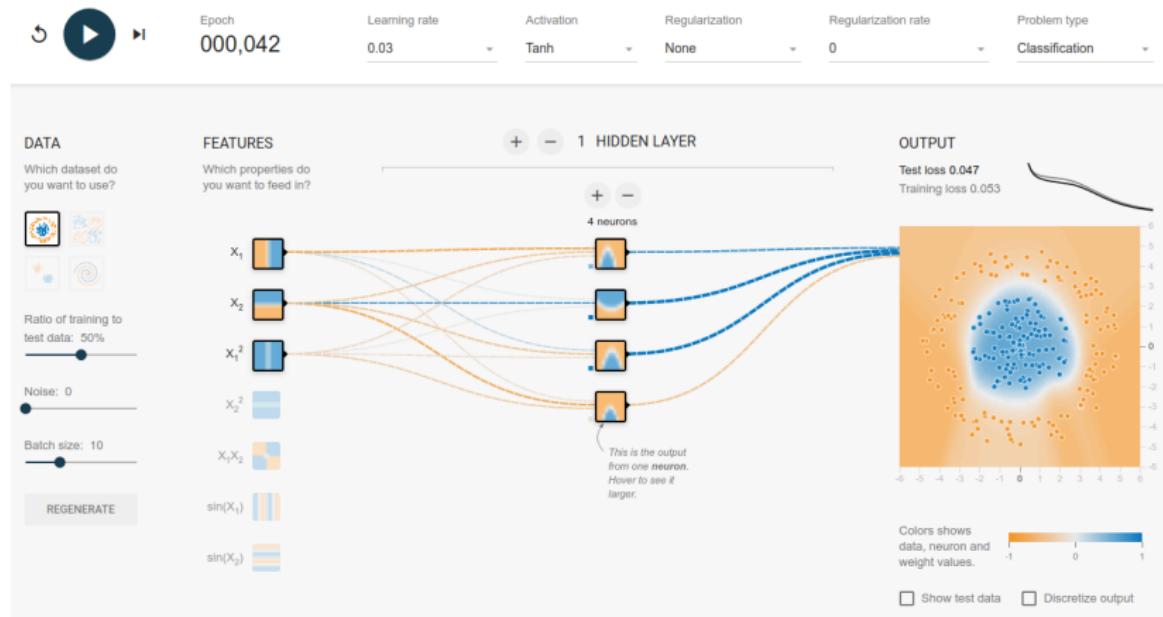


## a simple architecture (16 parameters)



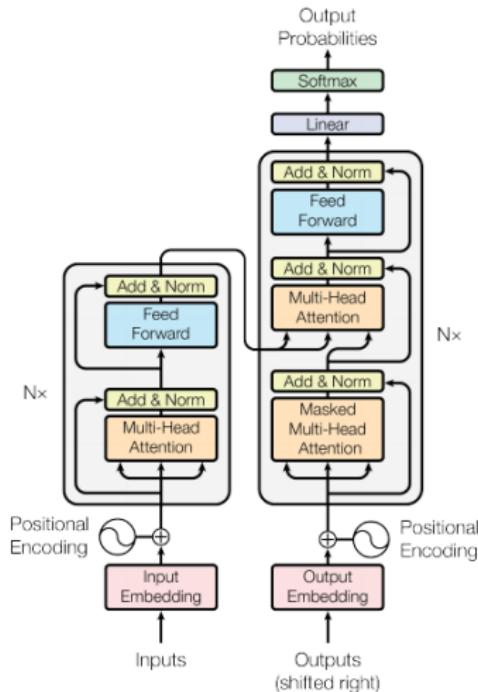
<https://playground.tensorflow.org>

## a simple architecture (16 parameters)

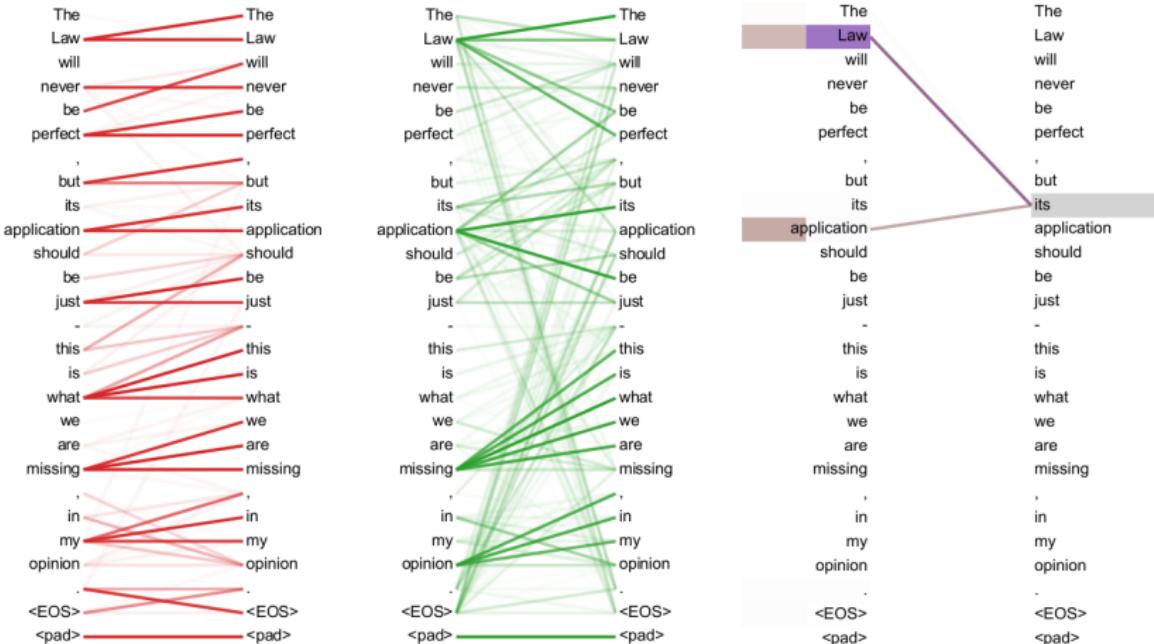


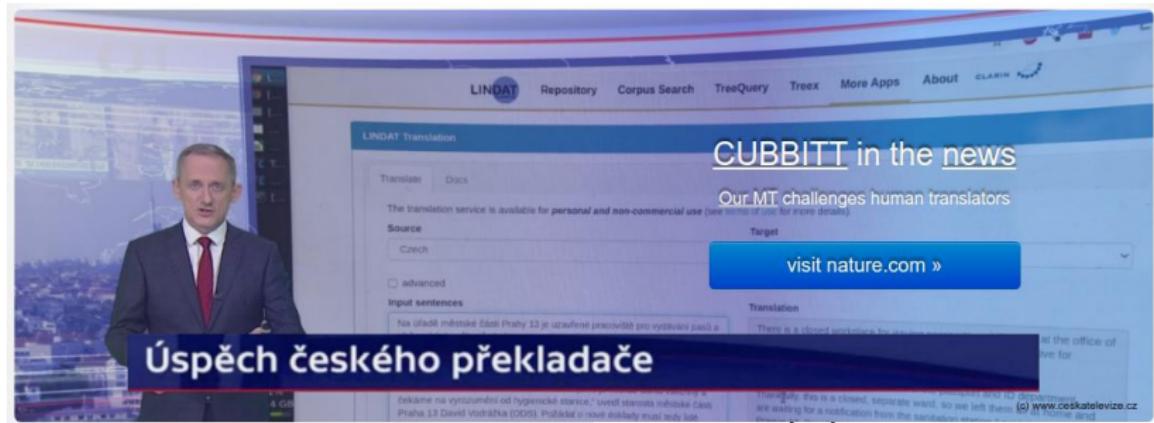
<https://playground.tensorflow.org>

## Transformer architecture (213 millions parameters)



the network learns important relationships between words  
 (in an unsupervised way, via self-attention layers)





Try CUBBITT at  
<https://lindat.cz/cubbitt>  
(En↔Cs, Fr, Pl)

nature communications

nature > nature communications > articles > article

Article | Open Access | Published: 01 September 2020

**Transforming machine translation: a deep learning system reaches news translation quality comparable to human professionals**

Martin Popel, Marketa Tomkova, Jakub Tomek, Łukasz Kaiser, Jakob Uszkoreit, Ondřej Bojar & Zdeněk Žabokrtsky

*Nature Communications* 11, Article number: 4381 (2020) | Cite this article

6273 Accesses | 76 Altmetric | Metrics

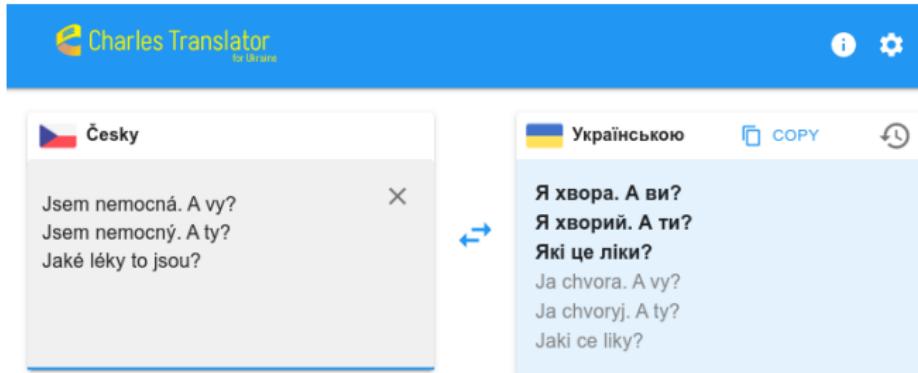


**Nový česko-ukrajinský překladač**

Zkuste si Charles Translator:

<https://translator.cuni.cz>





- první verze překladače vytvořena za 2 týdny
- cs→uk: asi 1400 uživatelů denně, uk-cs: asi 700
- přímý překlad (nikoli přes angličtinu)
- fonetický přepis azbuky, hlasový vstup,...
- API: překlad učebnic, PomahejUkrajine.cz,...

The screenshot shows the Charles Translator application for Ukraine. It has two main panels: one for Czech (left) and one for Ukrainian (right).  
Czech Panel (Left):  
- Flag: Czech flag (red, white, blue horizontal stripes)  
- Language: Česky  
- Text:

Jsem nemocná. A vy?  
Jsem nemocný. A ty?  
Jaké léky to jsou?

  
Ukrainian Panel (Right):  
- Flag: Ukrainian flag (blue and yellow horizontal stripes)  
- Language: Українською  
- Text:

Я хвора. А ви?  
Я хворий. А ти?  
Які це ліки?  
Ja chvora. A vy?  
Ja chvoryj. A ty?  
Jaki ce liky?

  
A blue double-headed arrow icon is positioned between the two panels.

The screenshot shows the Google Translate interface. It has two main panels: one for Czech (left) and one for Ukrainian (right).  
Czech Panel (Left):  
- Language: ROZPOZNAT JAZYK (Detektovat jazyk)  
- Language: ČEŠTINA  
- Text:

Jsem nemocná. A vy?  
Jsem nemocný. A ty?  
Jaké léky to jsou?

  
Ukrainian Panel (Right):  
- Language: UKRAJINŠTINA  
- Language: ČEŠTINA  
- Text:

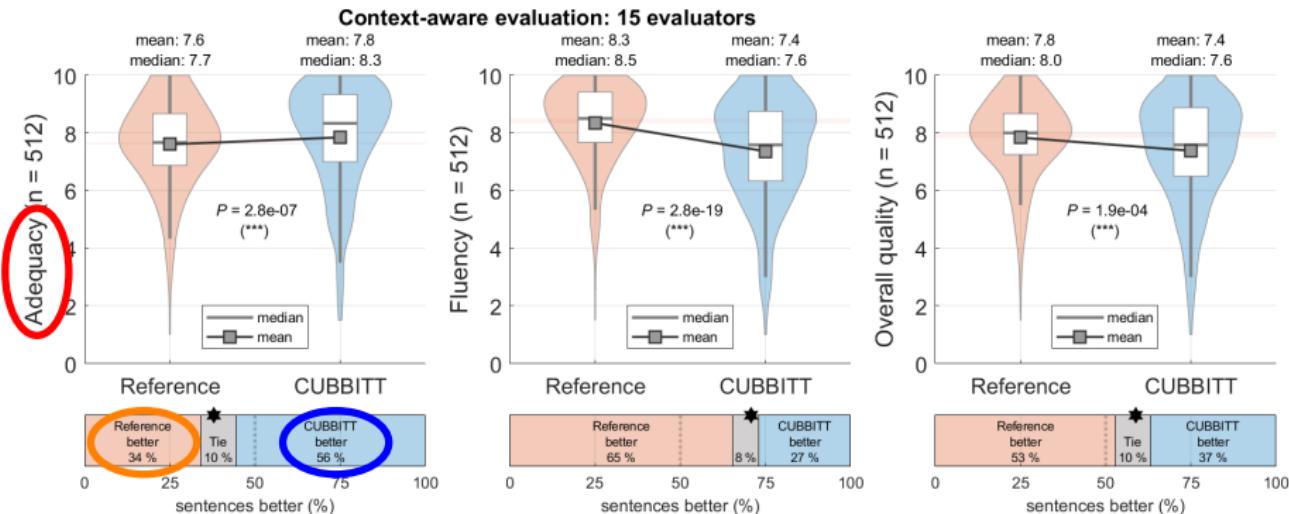
Я хворий. І ти?  
Я хворію. І ти?  
Які це наркотики?  
YA khvoryy. I ty?  
YA khvoriyu. I ty?  
Yaki tse narkotyky?

  
At the bottom, there are icons for microphone, speaker, and sharing, along with page navigation and a search bar.

# The main result of our human evaluation

16

56 % sentences translated more adequately by CUBBITT,  
34 % sentences by a professional translation agency



Support The Guardian    Subscribe    Find a job    Sign in

# The Guardian

News   Opinion   Sport   Culture   Lifestyle

World ► Europe US Americas Asia Australia Middle East Africa Inequality

[Facebook](#)

## Facebook translates 'good morning' into 'attack them', leading to arrest

Palestinian man questioned by Israeli police after embarrassing mistranslation of caption under photo of him leaning against bulldozer



zdroj: [The Guardian](#) (2017)

Support The Guardian    Subscribe    Find a job    Sign in

# The Guardian

News   Opinion   Sport   Culture   Lifestyle

World ► Europe US Americas Asia Australia Middle East Africa Inequality

Facebook

## Facebook translates 'good morning' into 'attack them', leading to arrest

Palestinian man questioned by Israeli police after embarrassing mistranslation of caption under photo of him leaning against bulldozer



zdroj: [The Guardian \(2017\)](#)

100+ billion words daily (Google, Microsoft, Baidu, Amazon,...)  
market size in 2020: USD 650 million, but rapidly growing

Source	Jana je žena. Pracuje jako průvodčí.
Google	Jana is a woman. He works as a guide.
Bing	Jana is a woman. He works as a conductor.
CUBBITT	Jane is a woman. He works as a conductor.
CUBBITT-doc	Jana is a woman. <b>She</b> works as a conductor.

Source	Saša je muž. Pracuje jako průvodčí.
DeepL	Sasha is a man. She works as a conductor.
CUBBITT-doc	Sasha is a man. <b>He</b> works as a conductor.

Source	Pavla není muž. Pracuje jako průvodčí.
DeepL	Paul is not a man. He works as a conductor.
CUBBITT-doc	Pavla is not a man. He works as a conductor.

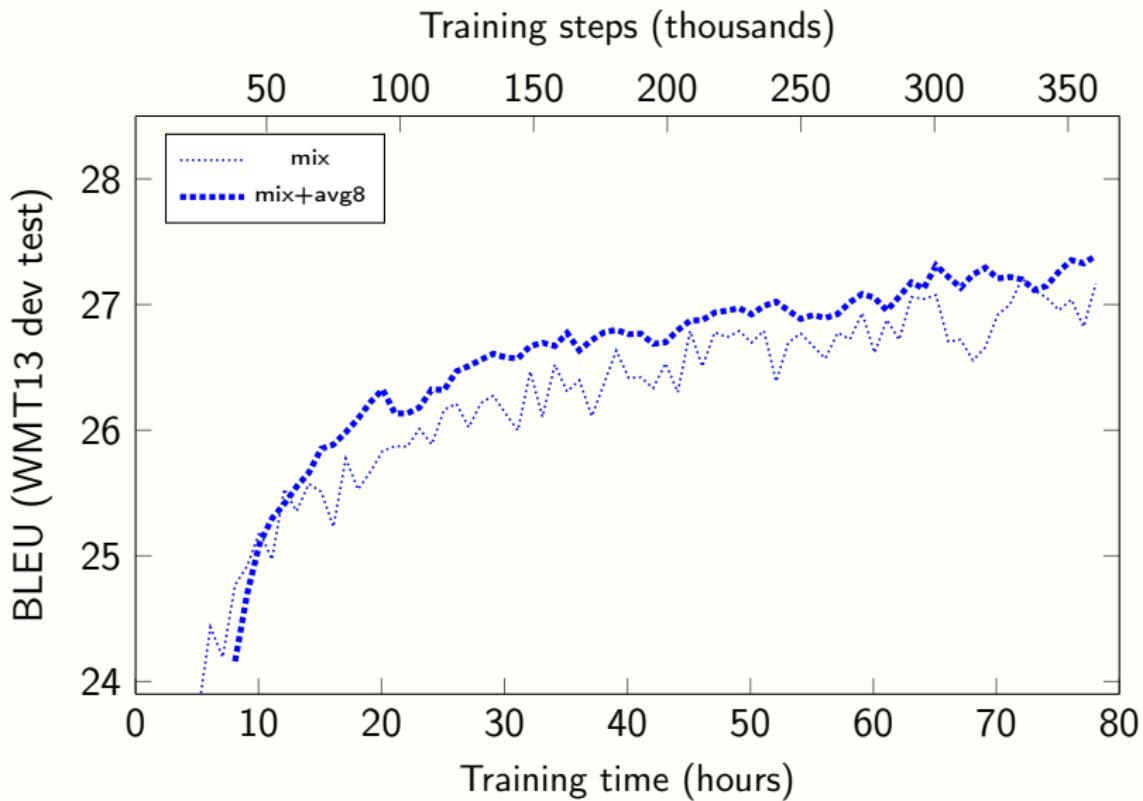
Want to know more?

- For EN→CS translation, we can exploit monolingual CS data.
- Translate the data back to English (with any CS→EN MT).
- Prepare synthetic parallel data (orig-CS, synth-EN).
- Train on both authentic and synthetic
  - **fine-tune BT**: first auth then auth+synth
  - **mix BT**: shuffle auth and synth sentences 1:1



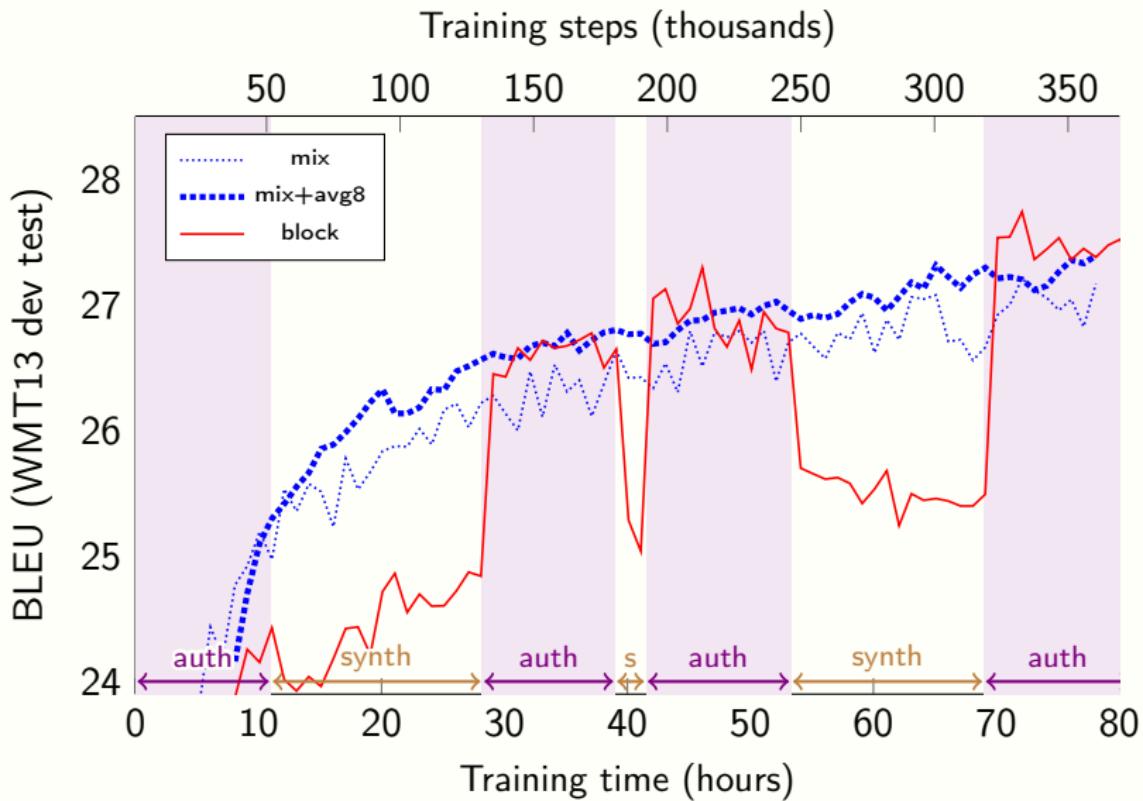
- For EN→CS translation, we can exploit monolingual CS data.
- Translate the data back to English (with any CS→EN MT).
- Prepare synthetic parallel data (orig-CS, synth-EN).
- Train on both authentic and synthetic
  - **fine-tune BT**: first auth then auth+synth
  - **mix BT**: shuffle auth and synth sentences 1:1
  - **block BT**: no shuffle, just concatenate auth and synth blocks





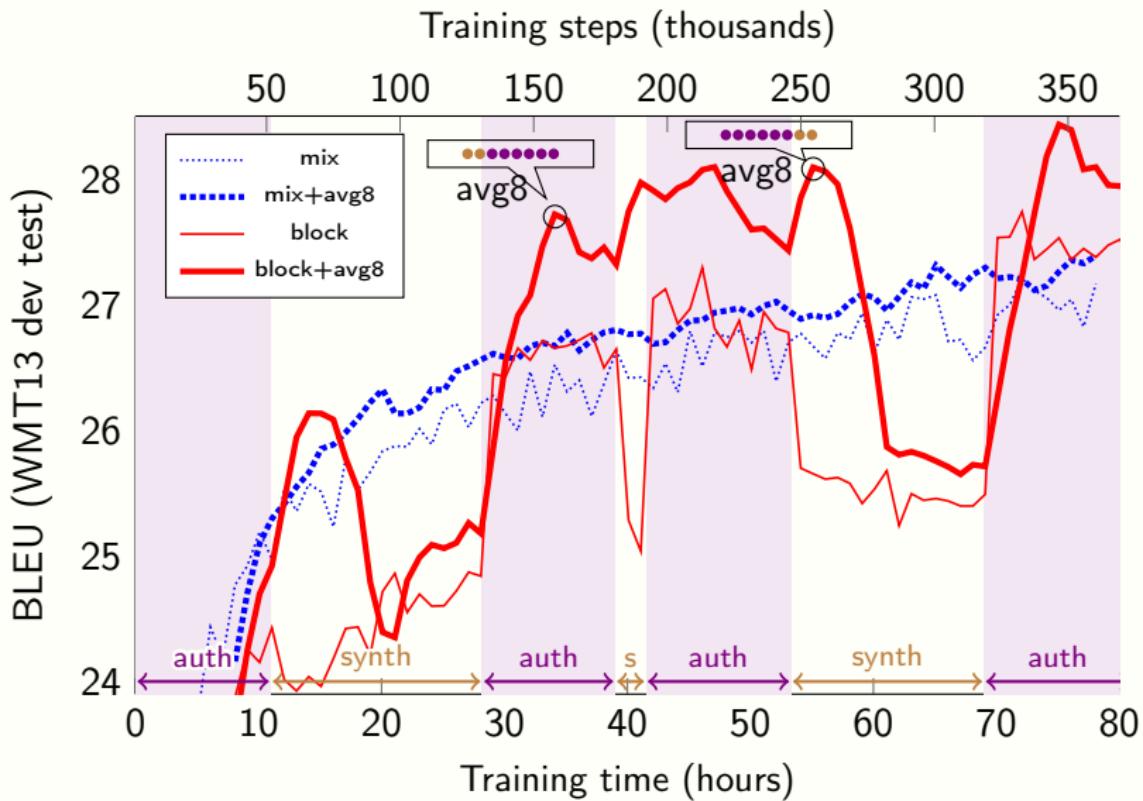
# Block Backtranslation

21



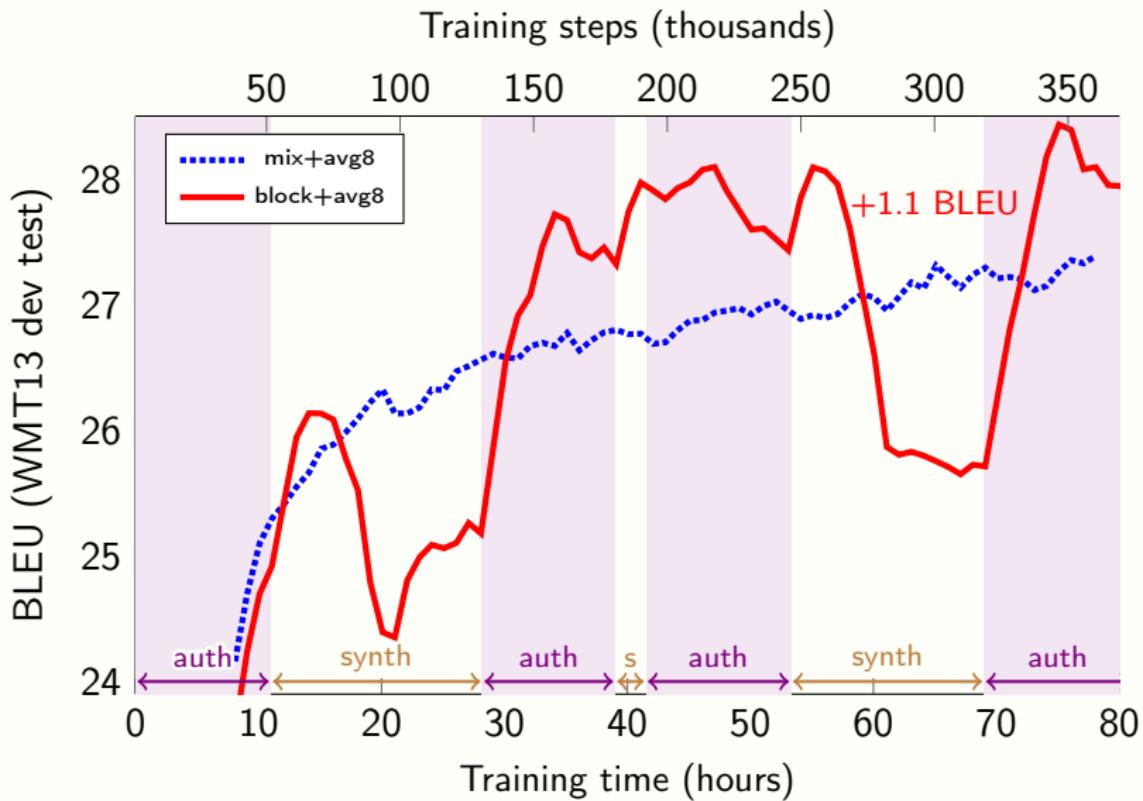
# Block Backtranslation

21



# Block Backtranslation

21



# Manual evaluation example

22

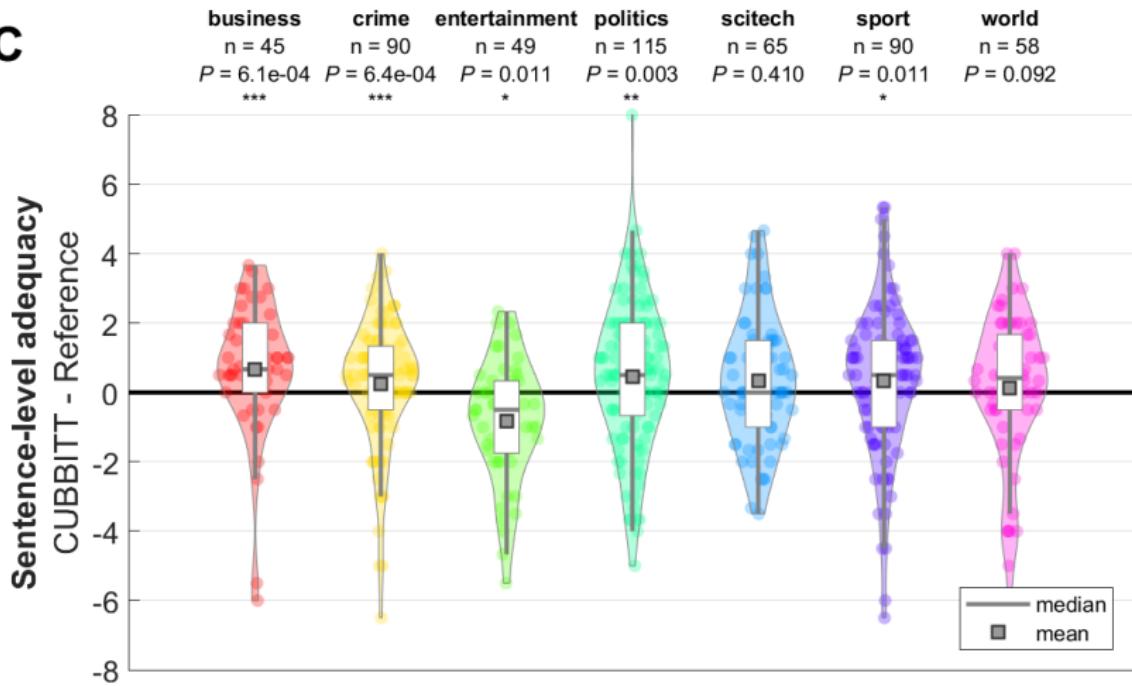
blind, sentence-level but document-aware, side-by-side (**RankME**)

	G	H	I	J	K	L	M	N	O	P
1	Source	Translation1	T1_overall	T1_adeqacy	T1_fluency	Translation2	T2_overall	T2_adeqacy	T2_fluency	Optional comment
168	"And we're protecting our shareholders from employment litigation."									
169	Companies started taking ethics, values and employee engagement more seriously in 2002 after accounting firm Arthur Andersen collapsed because of ethical violations from the Enron scandal, Quinal said.									
170	But it wasn't until "social media came into its own" that companies realized they couldn't stop their dirty laundry from going viral online.									
171	"Prior to using technology to monitor ethics, people used hope as a strategy," he said.									
172	Both Glint and Convergent offer their software as a service, charging companies recurring fees to use their products.									
173	It's a business model and opportunity that has the approval of venture capital investors, who have propped up both start-ups.	Je to obchodní model a příležitost, kterou schvaluji odvážní kapitáloví investoři, jenž podporují oba startupy.	7	6	7	Je to obchodní model a příležitost, která má souhlas investorů rizikového kapitálu, kteří podporují oba start-upy.	10	10	10	T1: chybný překlad termínu "venture capital"
174	Convergent raised \$10 million in funding in February from firms such as Sapphire Ventures and Tola Capital, bringing its total capital raised to \$47 million.	Convergent vybral v rámci své únorové kampaně od frenů jako Sapphire Ventures a Tola Capital celkově 10 milionů \$. A nakonec si odnesl kapitál ve výši 47 milionů \$.	3	4	3	Convergent ziskal v únoru finanční prostředky ve výši 10 milionů dolarů od firem jako Sapphire Ventures a Tola Capital, čímž se jeho celkový kapitál zvýšil na 47 milionů dolarů.	10	10	10	
175	Glint secured \$10 million in November from Bessemer Venture Partners, bringing its total funding to \$60 million.	Glint ziskal v listopadu 10 milionů \$ od Bessemer Venture Partners a v průběhu celé kampaně ziskal 60 milionů \$.	5	4	5	Glint ziskal v listopadu 10 milionů dolarů od společnosti Bessemer Venture Partners, čímž jeho celkové financování dosáhlo 60 milionů dolarů.	10	10	10	
176	These investments hardly come as a surprise, given the interconnected nature of companies, culture and venture capital.	Tyto investice jsou stále překvapující vzhledem k vzájemné povaze společnosti, kultury a rizikovému kapitálu.	3	4	3	Tyto investice nejsou vzhledem k propojenosti společnosti, kultury a rizikového kapitálu zádným překvapením.	10	10	10	
177	There's a growing body of research showing today's employees expect more from their workplaces than before.	Narůstající počet výzkumů jasně potvrzuje, že dnešní zaměstnanci očekávají od svého pracoviště více než kdy dříve.	5	5	5	Roste množství výzkumů, které ukazují, že dnešní zaměstnanci očekávají od svých pracovišť více než dříve.	10	10	10	
178	In competitive markets such as Silicon Valley, high salaries and interesting projects are merely table stakes.	A na konkurenčních trzích, jakým je např. Silicon Valley, jsou hlavní výhodou vysoké platy a zajímavé projekty.	6	5	8	Na konkurenčních trzích, jako je Silicon Valley, jsou vysoké platy a zajímavé projekty pouhými sázkami u stolu.	7	8	7	problém: význam terminu "table stakes"
179	Employees want to feel that they're accepted and valued and that they're giving their time to a company with a positive mission.	Zaměstnanci chtějí vnímat, že jsou přijímáni a oceňováni a že věnují svůj čas společnosti, která usiluje o pozitivní poslání.	9	9	9	Zaměstnanci chtějí mít pocit, že jsou přijímáni a ceněni a že věnují svůj čas společnosti s pozitivním posláním.	10	10	9	

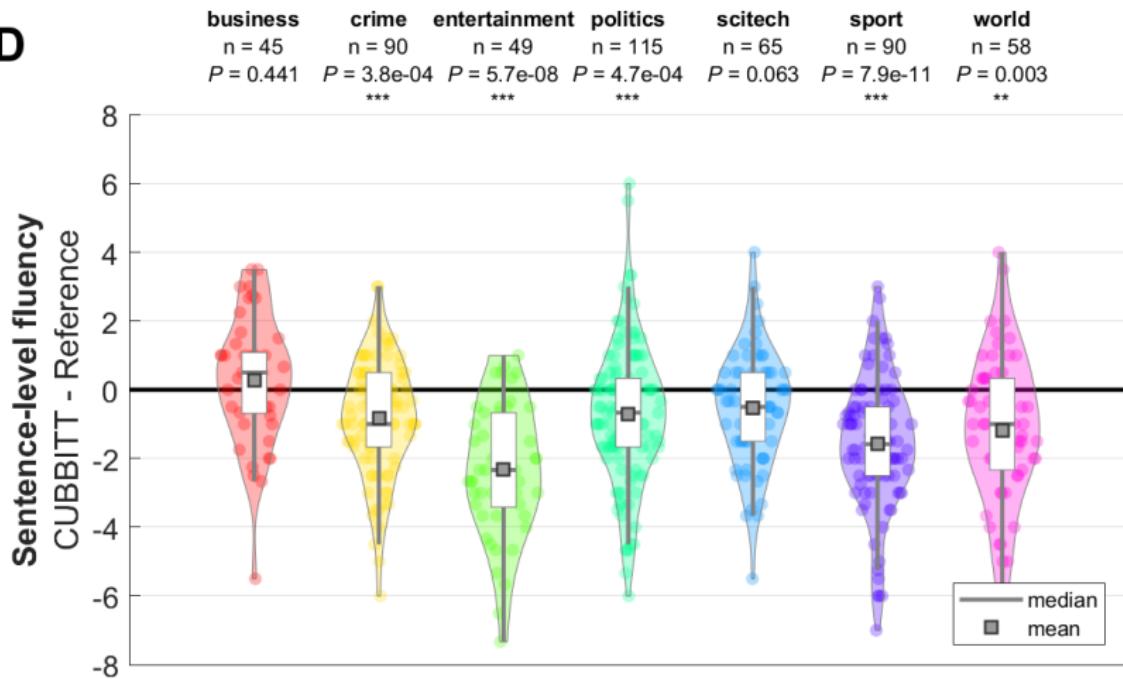
# It depends on the input text domain

23

C

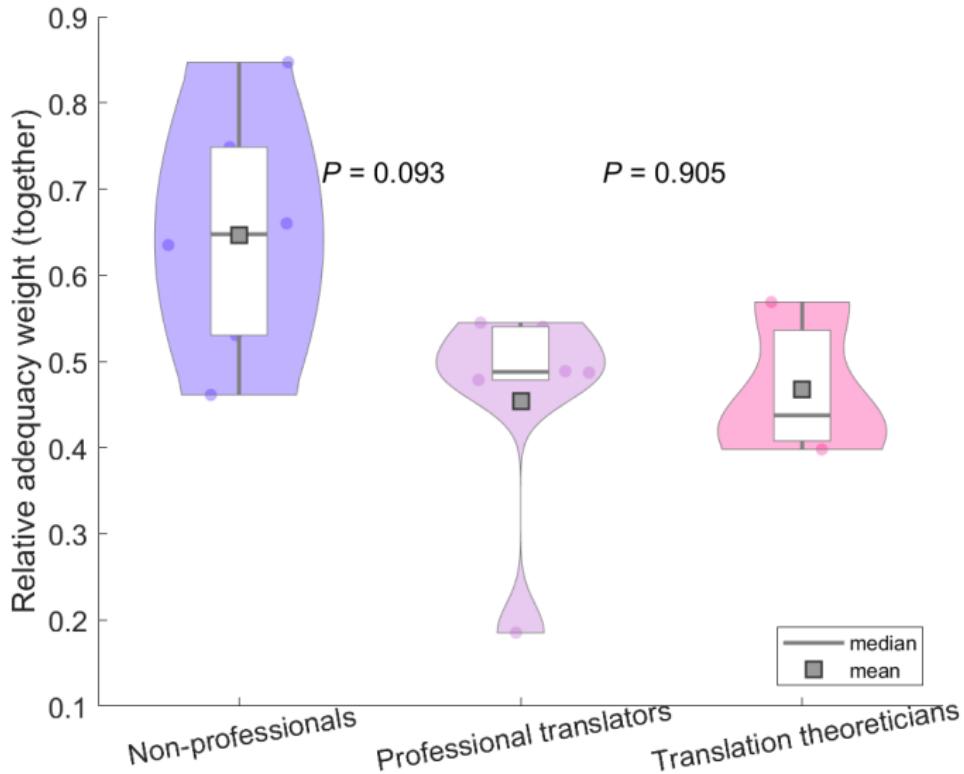


D

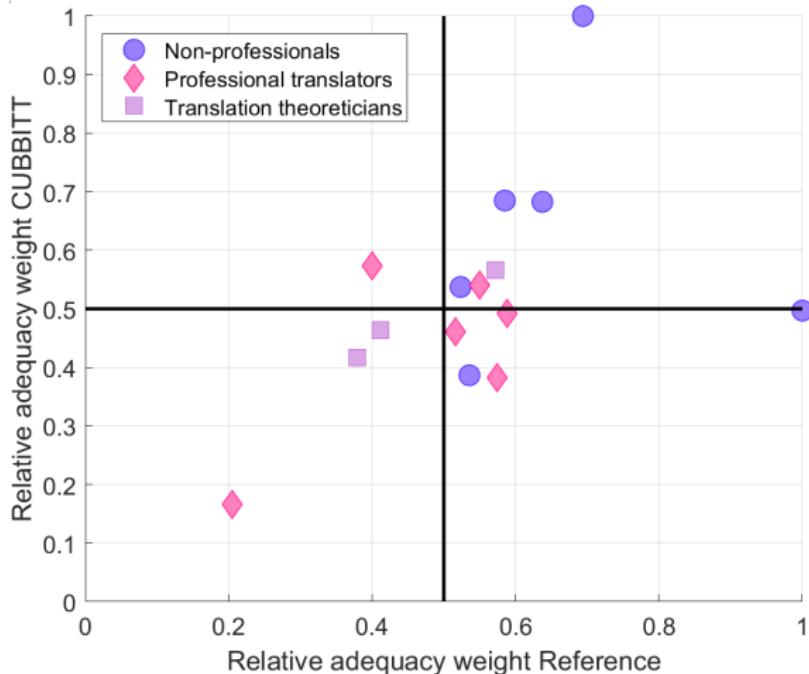


Can we predict the **overall quality**  
as a weighted average of **adequacy** and **fluency**?

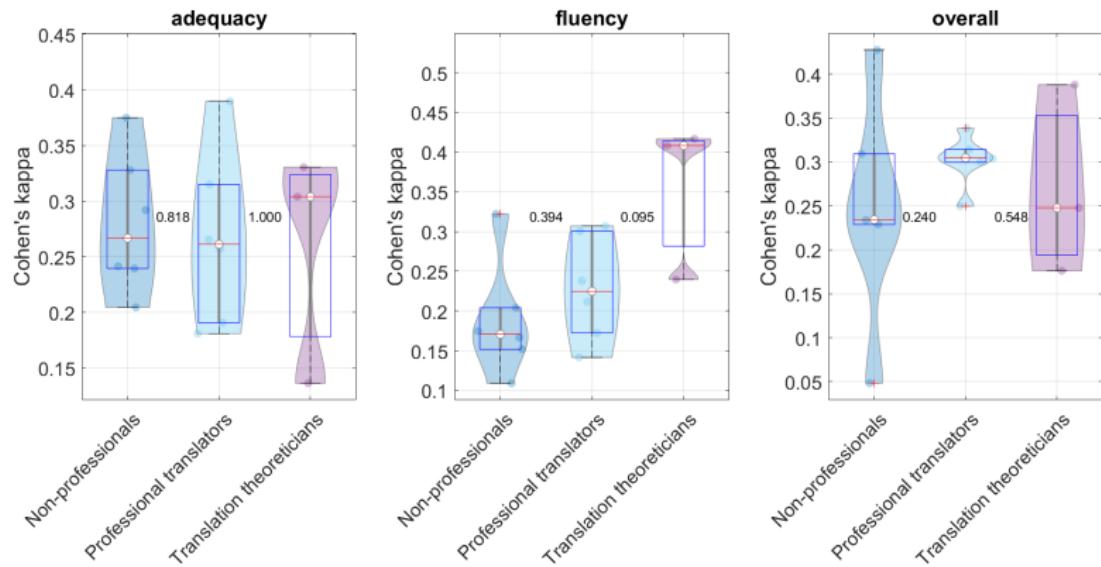
Translators put more emphasis on fluency,  
non-professionals on adequacy.



Some annotators put all the emphasis on adequacy,  
but only for one of the systems.

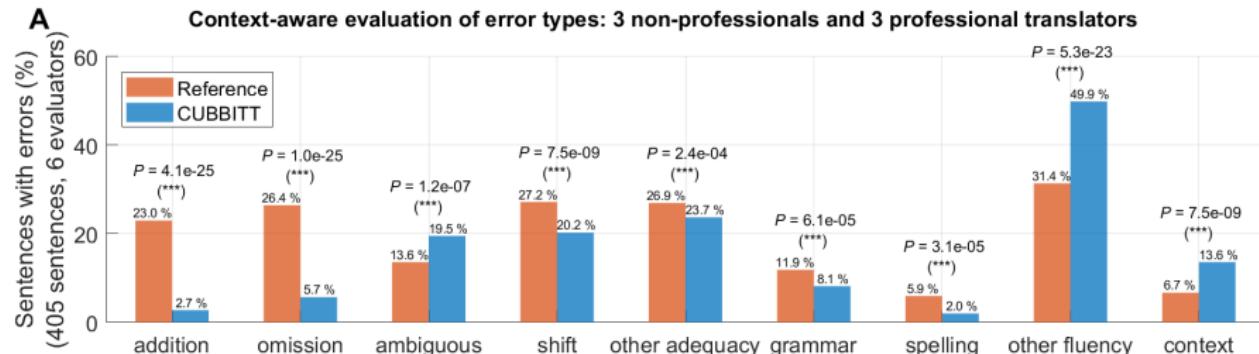


We cannot conclude that professionals are more reliable.



CUBBITT makes more errors than humans in

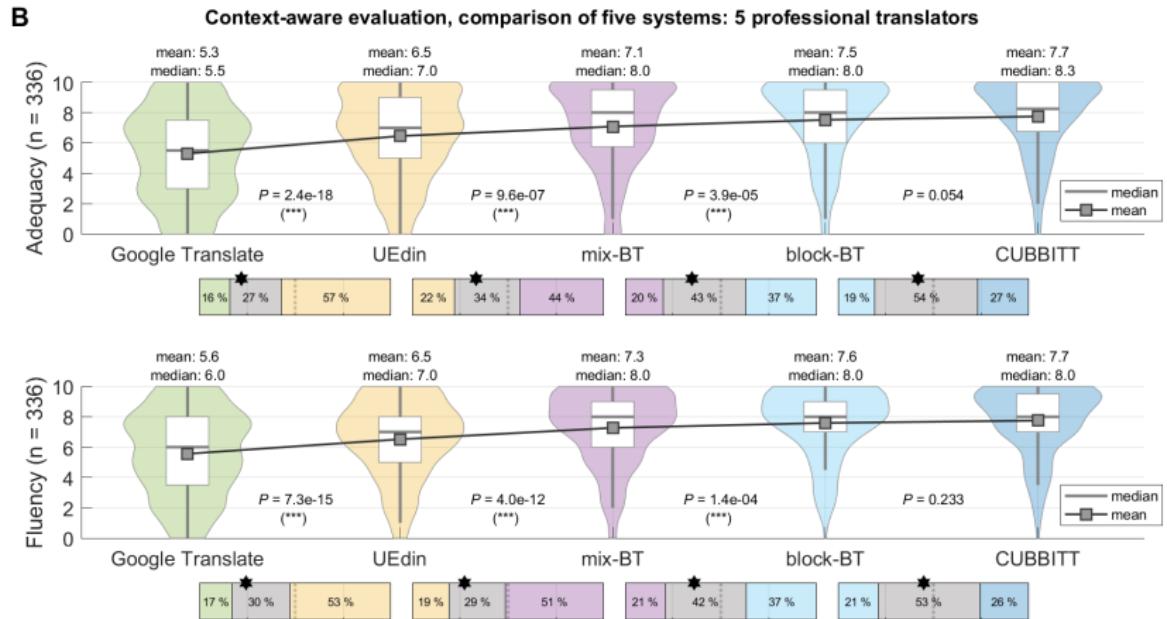
- translation of ambiguous words
- fluency (but not spelling and grammar) and
- cross-sentence context



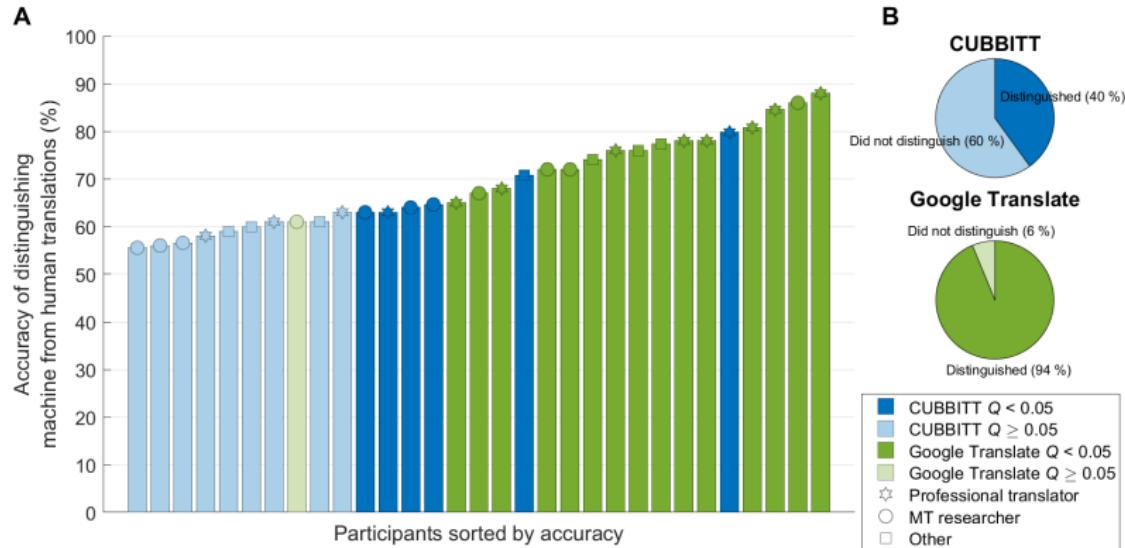
# Ablation analysis

29

BlockBT improves the quality (+0.4 adequacy, +0.3 fluency).



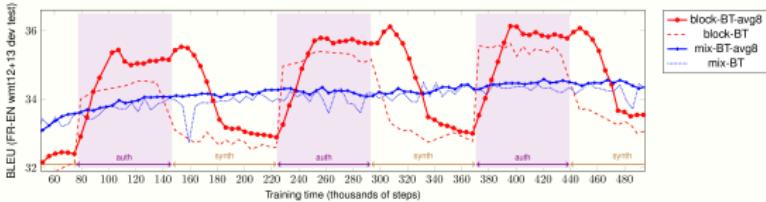
60 % of participants did not distinguish CUBBITT from human translations (on 100 isolated sentences)



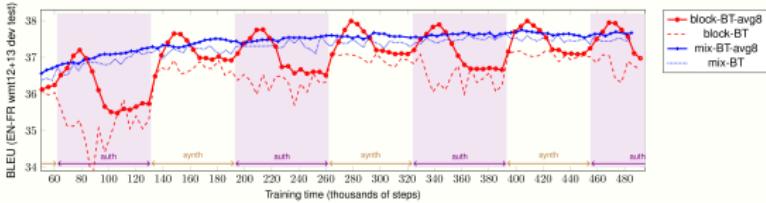
# BLEU Training curves for en-fr and en-pl

**A**

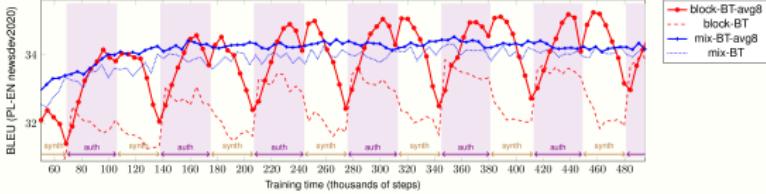
French - English

**B**

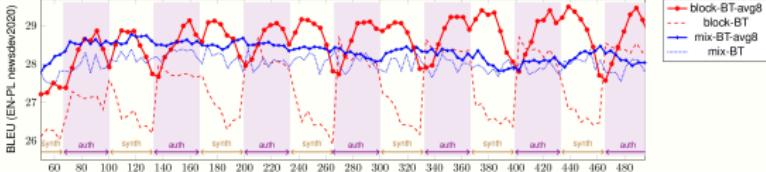
English - French

**C**

Polish - English

**D**

English - Polish



# Why Block-BT works?

32

**A**

Source sentence: "He was an original **guy** and lived life to the full" said **Gray** in a statement.



Translation of block-BT-Avg: "Byl to originální **chlap** a žil život naplno," uvedl **Gray** v prohlášení.

**Checkpoint 1 (SYNTH):** "Byl to originální **chlap** a žil život naplno" uvedl **Šedivý** v prohlášení.

**Checkpoint 2 (SYNTH):** "Byl to originální **chlap** a žil život naplno" uvedl **Šedivý** v prohlášení.

**Checkpoint 3 (AUTH):** „Byl to originální **člověk** a žil život naplno,“ uvedl **Gray** v prohlášení.

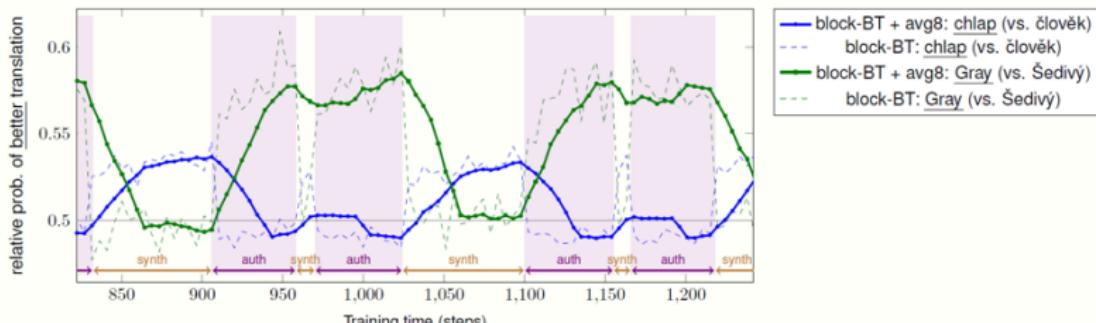
**Checkpoint 4 (AUTH):** „Byl to originální **člověk** a žil život naplno,“ uvedl **Gray** v prohlášení.

**Checkpoint 5 (AUTH):** „Byl to originální **člověk** a žil život naplno,“ uvedl **Gray** v prohlášení.

**Checkpoint 6 (AUTH):** „Byl to originální **člověk** a žil život naplno,“ uvedl **Gray** v prohlášení.

**Checkpoint 7 (AUTH):** „Byl to originální **člověk** a žil život naplno,“ uvedl **Gray** v prohlášení.

**Checkpoint 8 (AUTH):** „Byl to originální **člověk** a žil život naplno,“ uvedl **Gray** v prohlášení.

**B**

source	As good be an addled egg as an idle bird.
Bing	Jako dobrý být popolený vejce jako nečinný pták.
Google	Jako dobrá být včleněná vejce.
T2009	Dobré je fetácké vejce jako činný pták.
T2018	Dobří bud'te plete vejce jako nečinný pták.
CUBBITT	Stejně dobré je být pomateným vejcem jako zahálejícím ptákem.

source	A miss by an inch is a miss by a mile.
Bing	Miss o palec je Miss o míli.
Yandex	Slečna tím, že palec je vedle o míli.
Google	Chybějící palcem je míle vzdálená míle.
T2009	Slečna palec je slečna milionu.
T2018	Slečna palce je slečna míle.
CUBBITT	Minutí o centimetr je o kilometr.

Birds of a feather flock together.
Ptáci peří stáda dohromady.
Vrána k vráně sedá.
Vrána k vráně sedá.
Ptáci v bederním hejnu spolu.
Ptáci pérového hejna spolu.
Vrána k vráně sedá.

Try Charles Translator (CUBBITT) at:  
<https://translator.cuni.cz>