

Esther Ploeger

Department of Computer Science Aalborg University, Copenhagen

espl@cs.aau.dk



The start of a research project...



The start of a research project...





Some facts:

The official language of Greenland



Some facts:

- The official language of Greenland
- Family: Inuit-Yupik-Unangan



Some facts:

- The official language of Greenland
- Family: Inuit-Yupik-Unangan
- Speakers: approx. 50,000 60,000



Some facts:

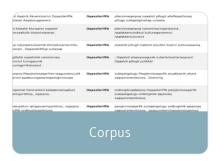
- The official language of Greenland
- Family: Inuit-Yupik-Unangan
- Speakers: approx. 50,000 60,000
- Greenland was under Danish colonial rule from 1721 until 1953











[ɔ ²qa: si lɜ ¹aif ¹fim ²mi:# ɔ ²qa: ¹sɜp ¹pas su ^wa lɜ ¹aif fik# ¹im ¹mik kut# su li ^jɑ ¹qɑf fi ^ju ^wɔq#]

IPA Converter











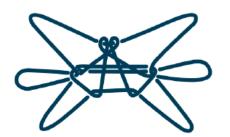
What could we work on?







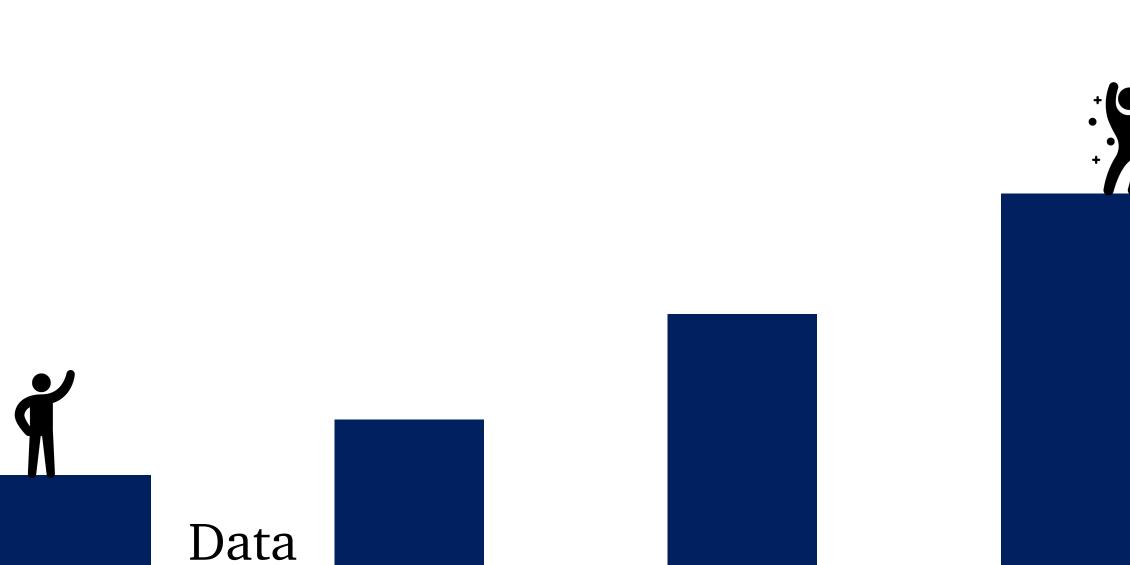
"An area in need of improvement is danish-greenlandic law texts."



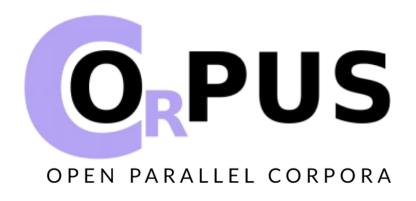








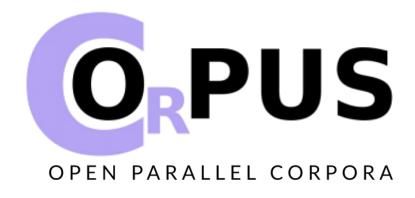




Tiedemann (2009)

Size: 291 lines

Domain: manuals



Tiedemann (2009)

Size: 291 lines

Domain: manuals



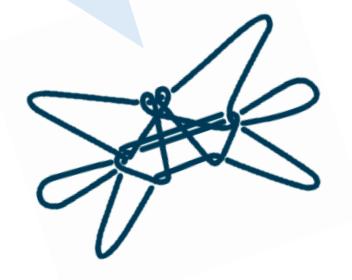
Kristensen-Mclachlan and Nedergard (2024)

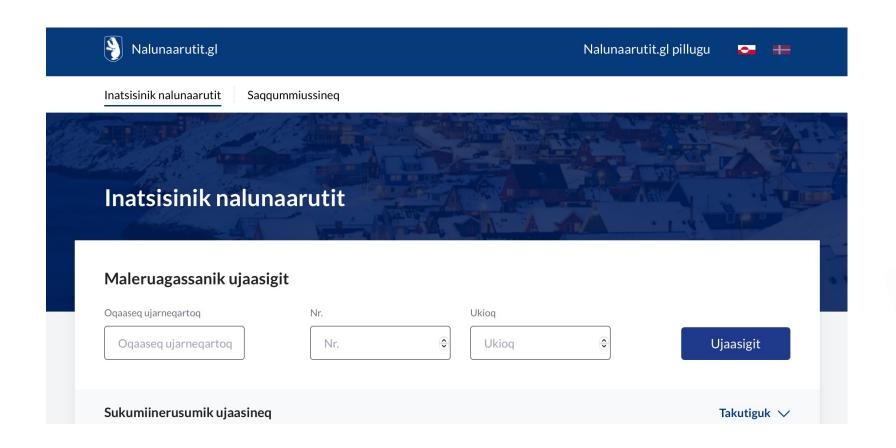
Size: 1.2m West-Greenlandic

words; 2.1m Danish

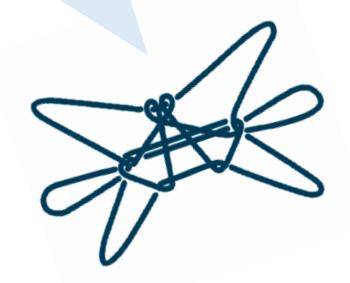
Domain: news

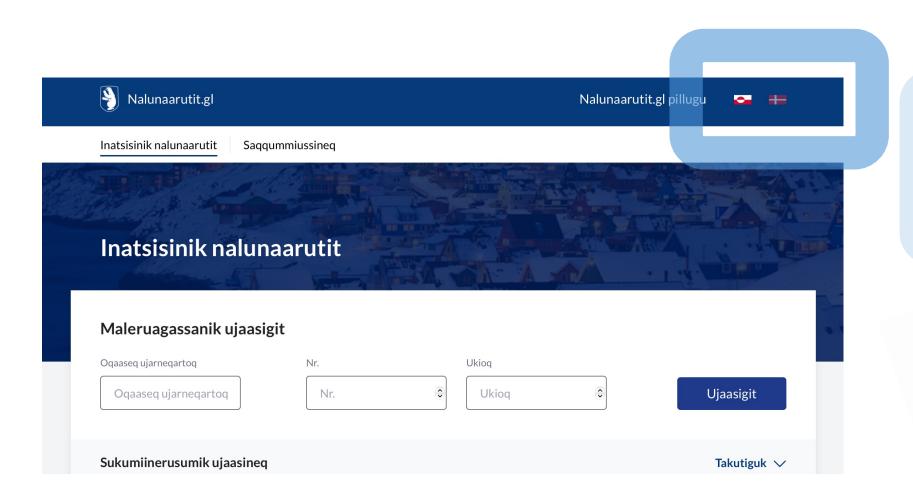
There is publicly available legal text!



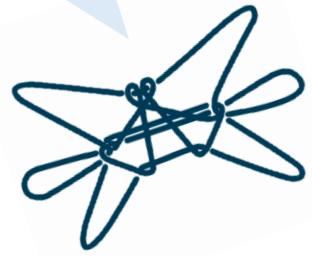


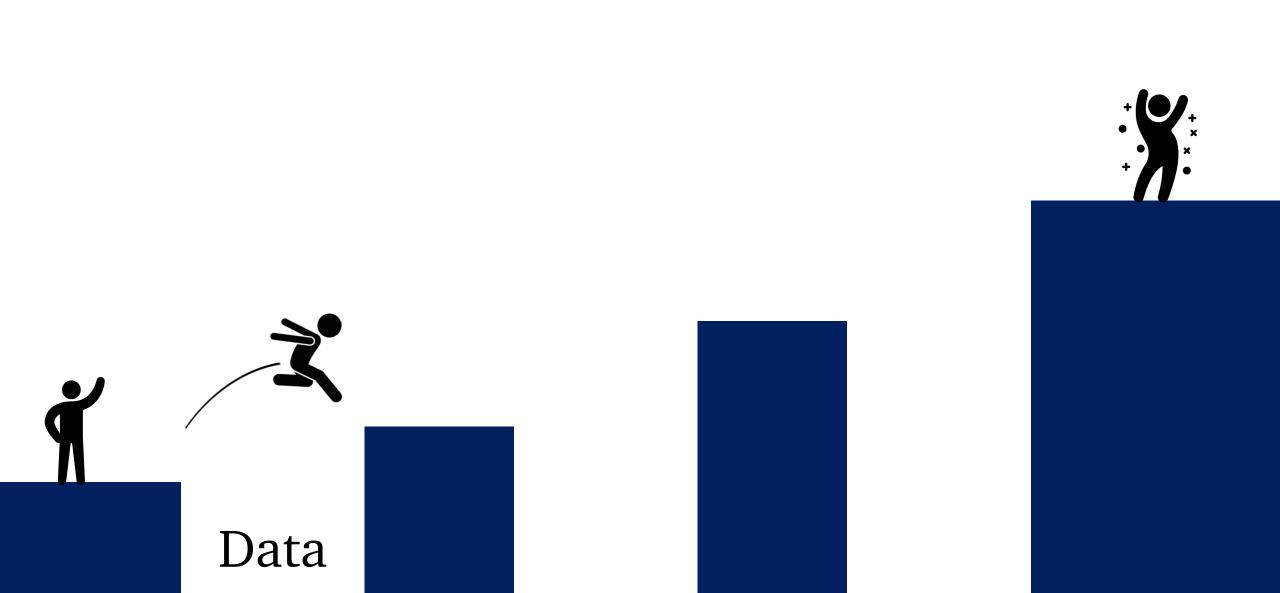
There is publicly available legal text!

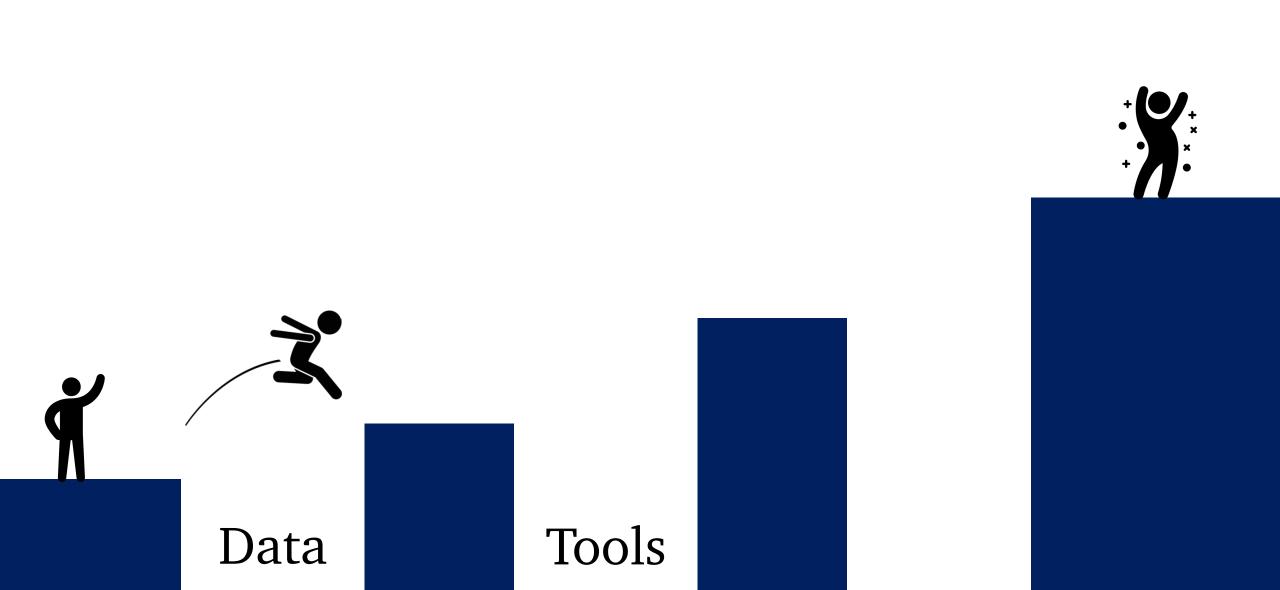




There is publicly available legal text!







Sentence Alignment

Sentence Alignment

e1: Bob isn't here; he went home. f1: Bob n'est pas là.

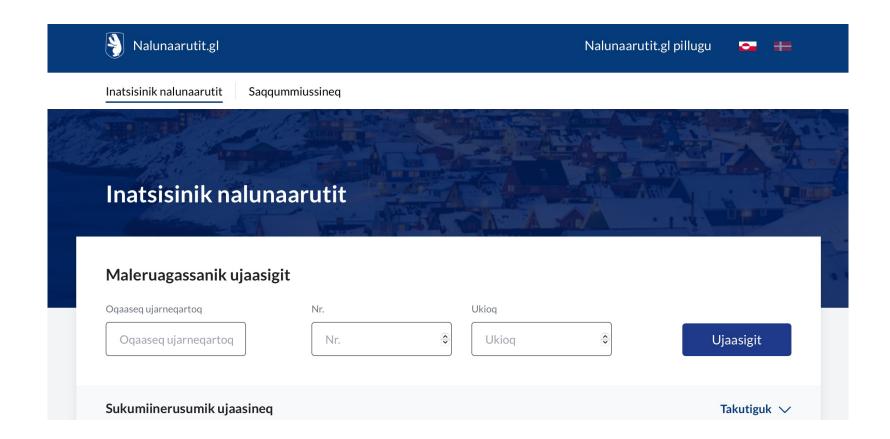
e2: Bob lives on the beach _____ f2: Bob est rentré chez lui.

e3: Bob has two cats. f3: Il vit sur la plage.

e4: Bob has a dog. f4: Bob a un chien et deux chats.

e5: Bob will be back tomorrow. \ f5: Le chat de Bob s'appelle Fluffy.

f6: Bob sera de retour demain.



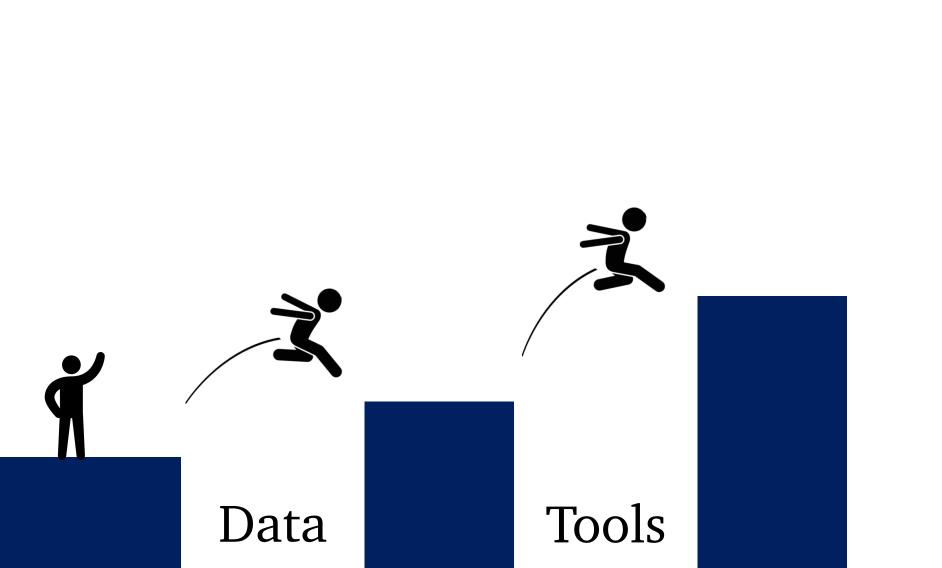
- § 1. Inatsit manna suliffeqarfinnut inunnullu makkununnga atuuppoq:
- 1) Suliffeqarfiit inuillu nunat allat aningaasaannut nuutsitsisarnermik inuussutissarsiuteqartut.
- 2) Suliffeqarfiit inuillu attartortitsinernut aningaasalersuinermik imaluunniit taarsigassarsisitsisarnermik inuussutissarsiuteqartut.
- 3) Illunik nioqquteqarnermi akunnermiliuttartut.
- 4) Suliffeqarfiit inuillu illunik nioqquteqarnermi akunnermiliuttartutut kiffartuussinernik inuussutissarsiuteqartut.
- § 1. Denne lov finder anvendelse på følgende virksomheder og personer:
- 1) Virksomheder og personer, der erhvervsmæssigt udøver virksomhed med valutaveksling.
- 2) Virksomheder og personer, der erhvervsmæssigt udøver finansiel leasing eller udøver udlånsvirksomhed.
- 3) Ejendomsmæglere.
- 4) Virksomheder og personer, der i øvrigt erhvervsmæssigt leverer samme ydelser som ejendomsmæglere.

- § 1. Inatsit manna suliffeqarfinnut inunnullu makkununnga atuuppoq:
- 1) Suliffeqarfiit inuillu nunat allat aningaasaannut nuutsitsisarnermik inuussutissarsiuteqartut.
- 2) Suliffeqarfiit inuillu attartortitsinernut aningaasalersuinermik imaluunniit taarsigassarsisitsisarnermik inuussutissarsiuteqartut.
- 3) Illunik nioqquteqarnermi akunnermiliuttartut.
- **4)** Suliffeqarfiit inuillu illunik nioqquteqarnermi akunnermiliuttartutut kiffartuussinernik inuussutissarsiuteqartut.
- § 1. Denne lov finder anvendelse på følgende virksomheder og personer:
- 1) Virksomheder og personer, der erhvervsmæssigt udøver virksomhed med valutaveksling.
- 2) Virksomheder og personer, der erhvervsmæssigt udøver finansiel leasing eller udøver udlånsvirksomhed.
- **3)** Ejendomsmæglere.
- **4)** Virksomheder og personer, der i øvrigt erhvervsmæssigt leverer samme ydelser som ejendomsmæglere.

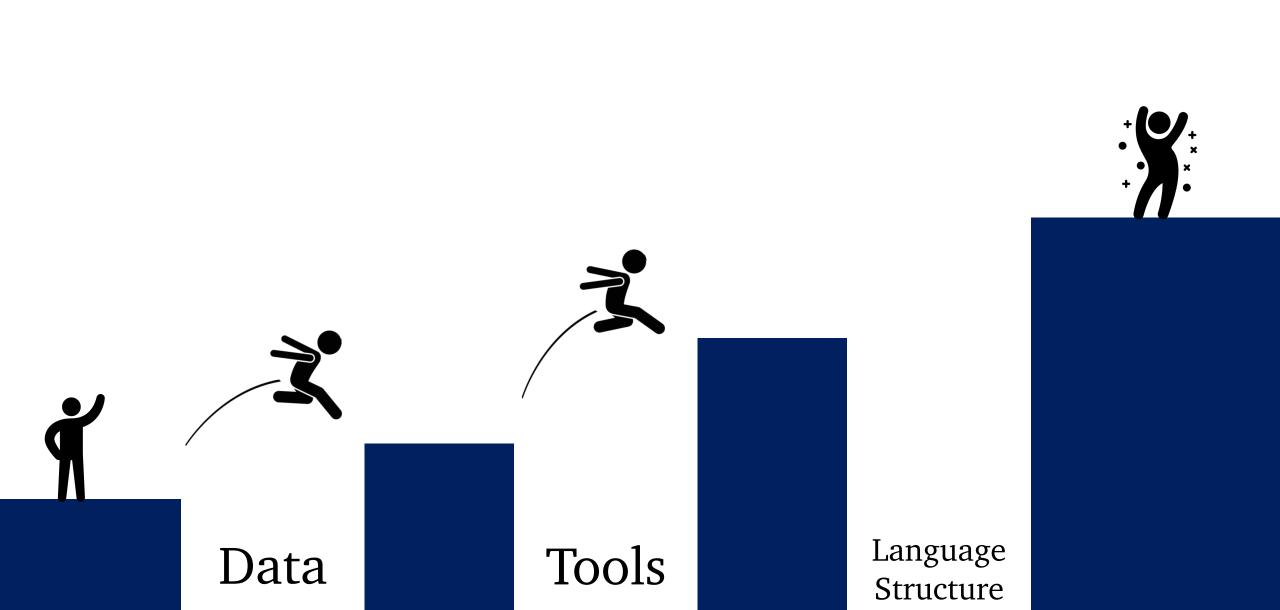
- § 1. Inatsit manna suliffeqarfinnut inunnullu makkununnga atuuppoq:
- 1) Suliffeqarfiit inuillu nunat allat aningaasaannut nuutsitsisarnermik inuussutissarsiuteqartut.
- 2) Suliffeqarfiit inuillu attartortitsinernut aningaasalersuinermik imaluunniit taarsigassarsisitsisarnermik inuussutissarsiuteqartut.
- 3) Illunik nioqquteqarnermi akunnermiliuttartut.
- **4)** Suliffeqarfiit inuillu illunik nioqquteqarnermi akunnermiliuttartutut kiffartuussinernik inuussutissarsiuteqartut.
- § 1. Denne lov finder anvendelse på følgende virksomheder og personer:
- 1) Virksomheder og personer, der erhvervsmæssigt udøver virksomhed med valutaveksling.
- **2)** Virksomheder og personer, der erhvervsmæssigt udøver finansiel leasing eller udøver udlånsvirksomhed.
- **3)** Ejendomsmæglere.
- **4)** Virksomheder og personer, der i øvrigt erhvervsmæssigt leverer samme ydelser som ejendomsmæglere.

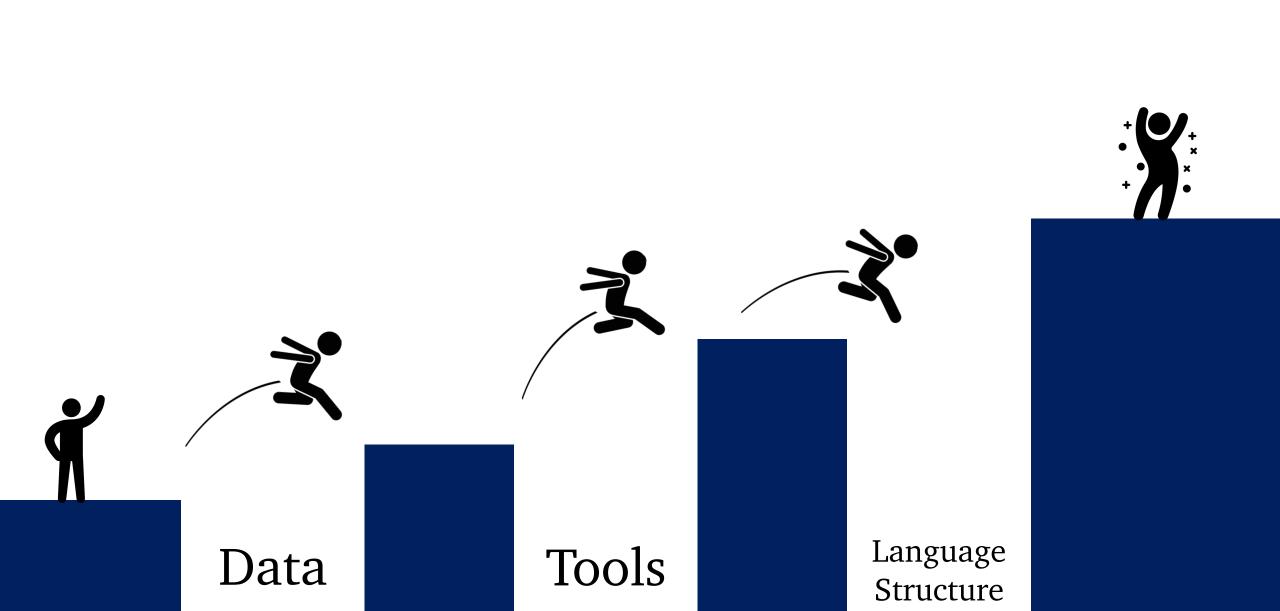


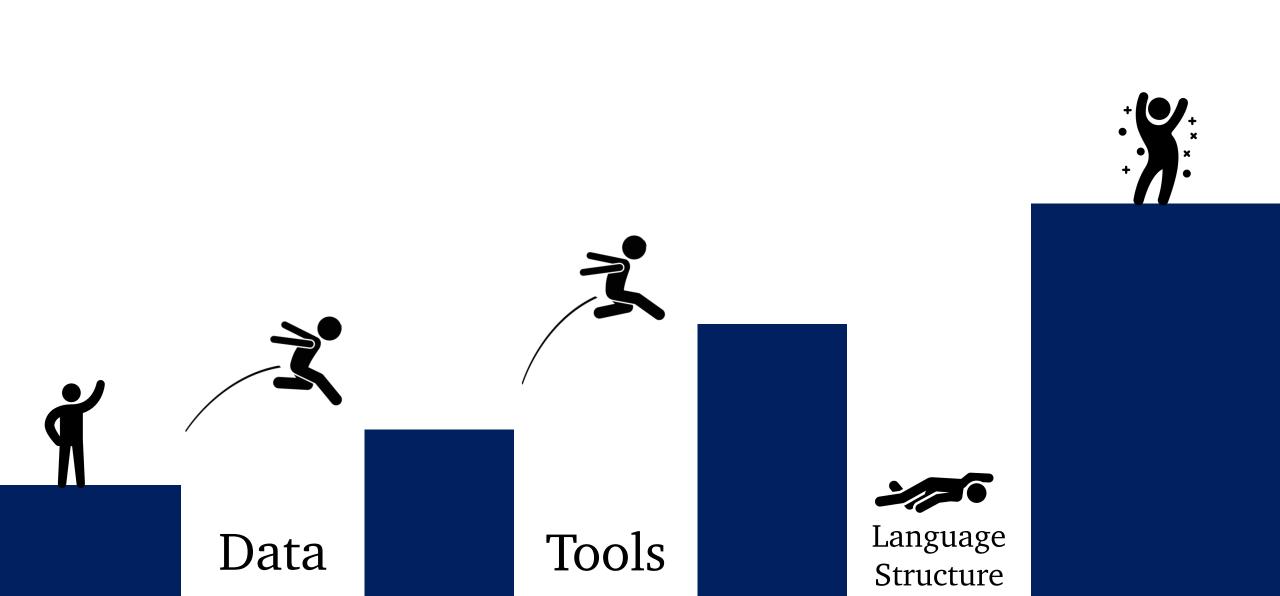
	# Lines		# Words	
Split	GL	DA	GL	DA
Training	39,936	39,936	663,734	929,904
Validation	1,000	1,000	16,594	23,021
Testing	1,000	1,000	16,665	23,846
Total	41,936	41,936	696,993	976,771











Obstacle 3: Language Structure

Obstacle 3: Language Structure

English: if he doesn't have a dog

Danish: hvis han ikke har en hund

Obstacle 3: Language Structure

English:

Danish:

Obstacle 3: Language Structure

English:

Danish:

if he doesn't have a dog

hvis han ikke har en hund

West-Greenlandic:

qimmeqanngikkuni

Obstacle 3: Language Structure

English:

Danish:

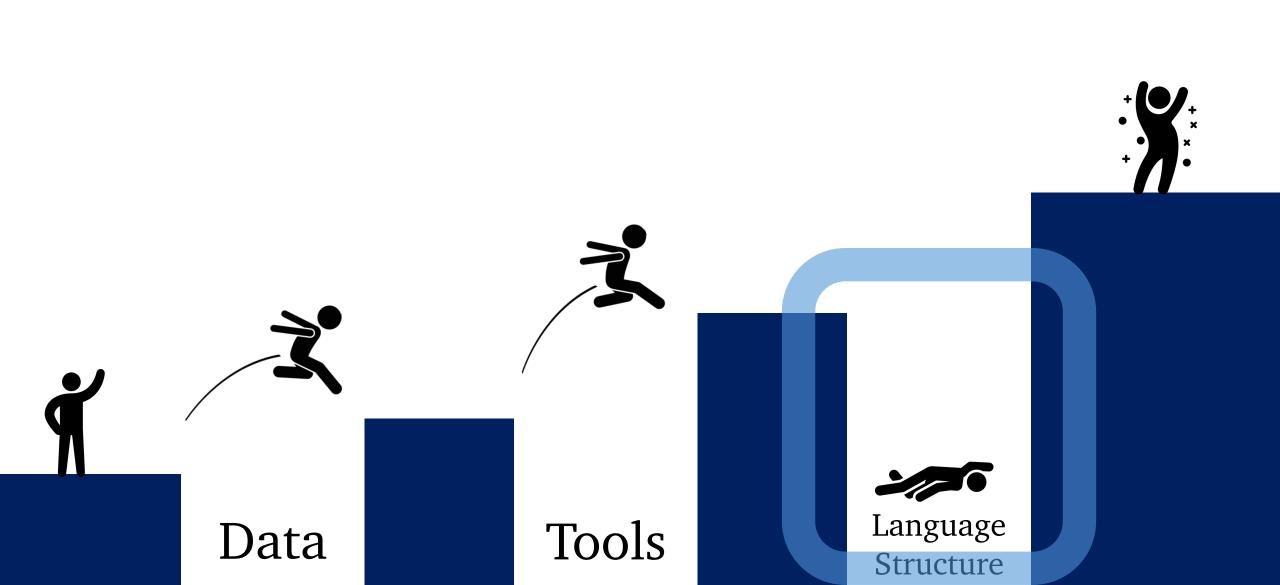
West-Greenlandic:

if he doesn't have a dog

hvis han ikke har en hund

qimmeqanngikkuni

A Tale of Three Obstacles



How should we segment text in a morphologically complex language like West-Greenlandic?

Why Segment Text?

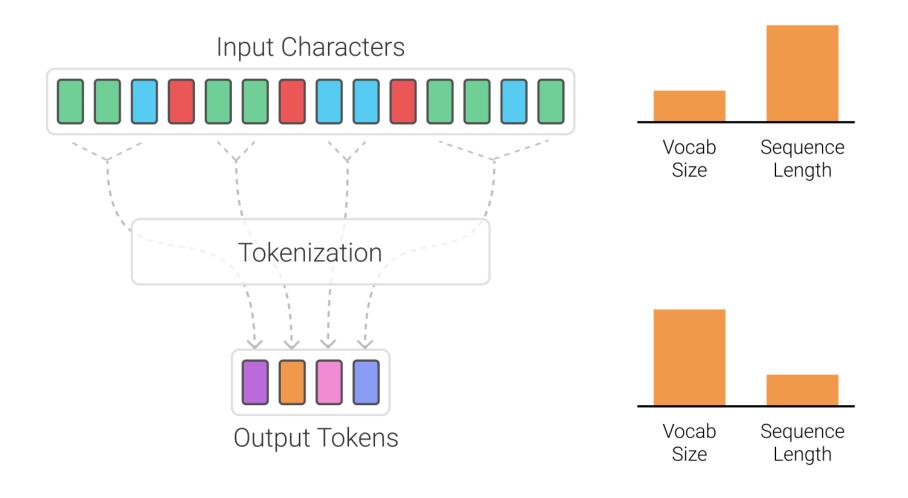


Image: https://www.thoughtvector.io/blog/subword-tokenization/

BPE

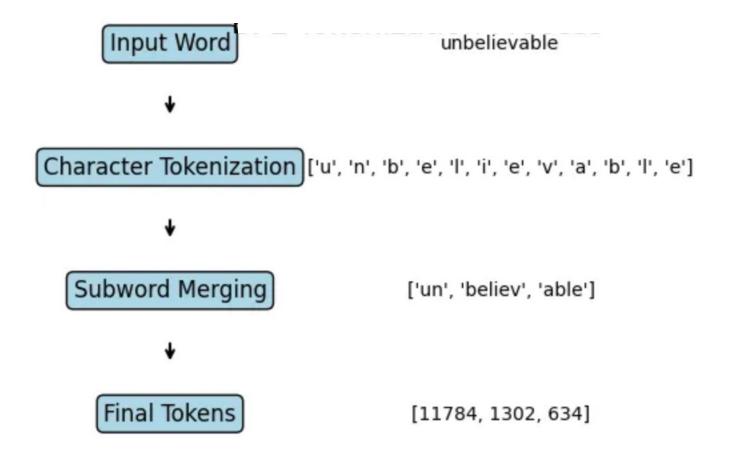
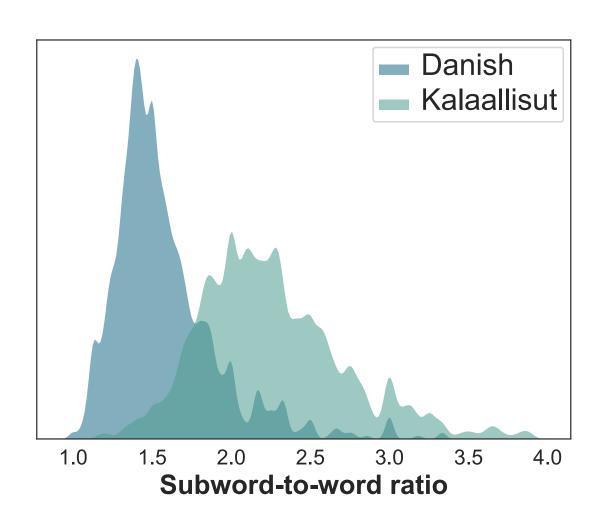


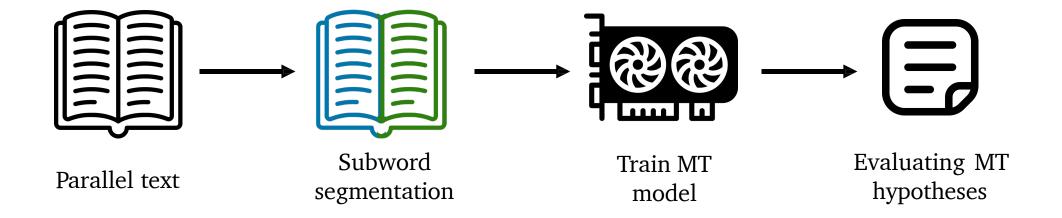
Image: https://python.plainenglish.io/efficient-tokenization-with-byte-pair-encoding-bpe-for-neural-networks-7cf4a54b5fd0

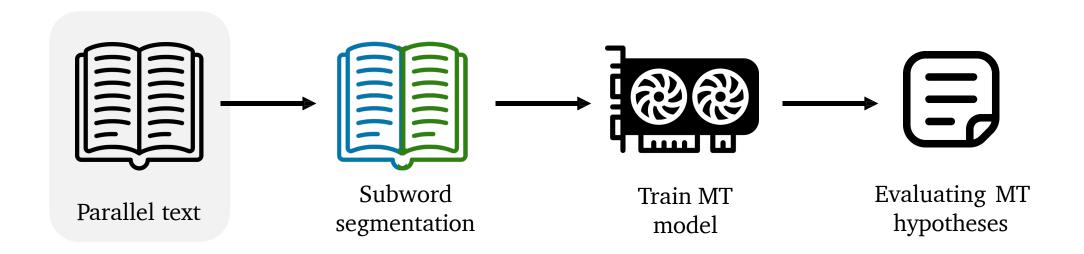
BPE

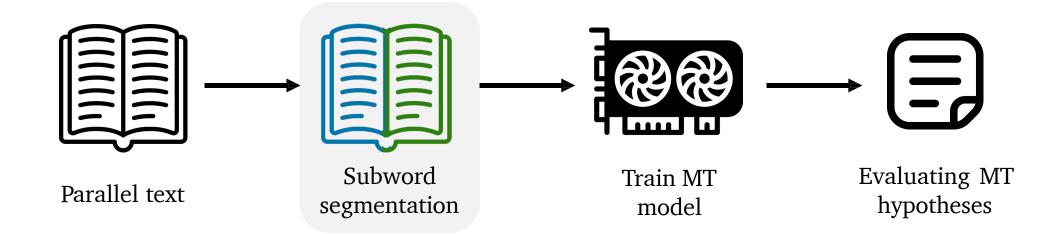
Biased by language structure!

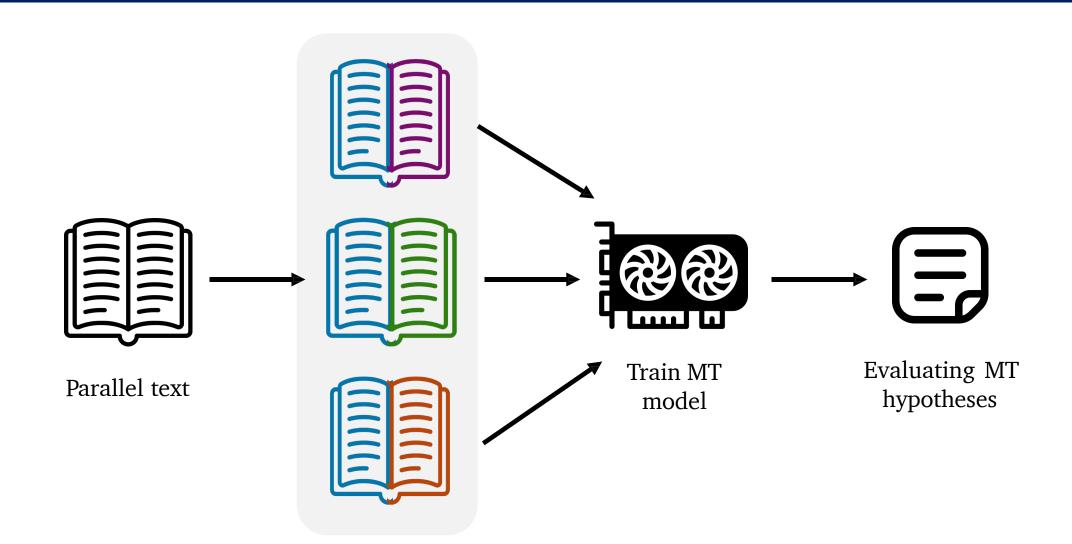


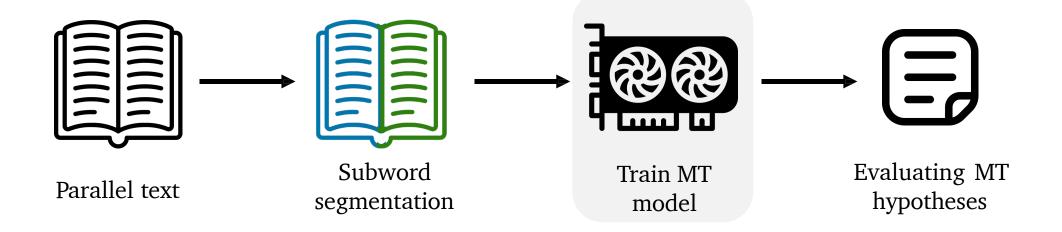
Previous work on neural polysynthetic language modelling suggests that BPE is suboptimal!

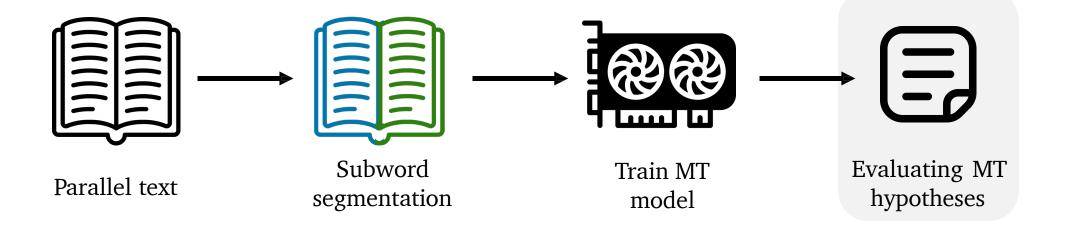












if he doesn't have a dog

hvis han ikke har en hund

qimmeqanngikkuni

if he doesn't have a dog

hvis han ikke har en hund

qimmeqanngikkuni

qimme qa nngik ku ni



qimme qa nngik ku ni

GPT-4: qimmeqanngikkuni

$$P = \frac{ \left| \left\{ gold\ morphemes \right\} \ \cap \ \left\{ subwords \right\} \right| }{ \left| \left\{ subwords \right\} \right| }$$

$$R = \frac{ \mid \{ gold\ morphemes \} \ \cap \ \{ subwords \} \mid }{ \mid \{ gold\ morphemes \} \mid }$$

Method	Prec. (%)	Rec. (%)	F1 (%)
BPE	10.81	12.42	11.56
Unigram	30.88	<u>37.68</u>	<u>33.94</u>
Morfessor	<u>31.08</u>	31.61	31.34
FlatCat	29.58	29.40	29.49

Does this actually help machine translation?

Results: Machine Translation

	Machine Translation			
	Danish→Kalaallisut		Kalaallisut→Danish	
Method	chrF2	BLEU	chrF2	BLEU
None	44.6	3.4	61.5	15.4
BPE	56.4	8.7	<u>64.2</u>	<u>21.5</u>
Unigram	61.4	<u>10.1</u>	58.9	17.1

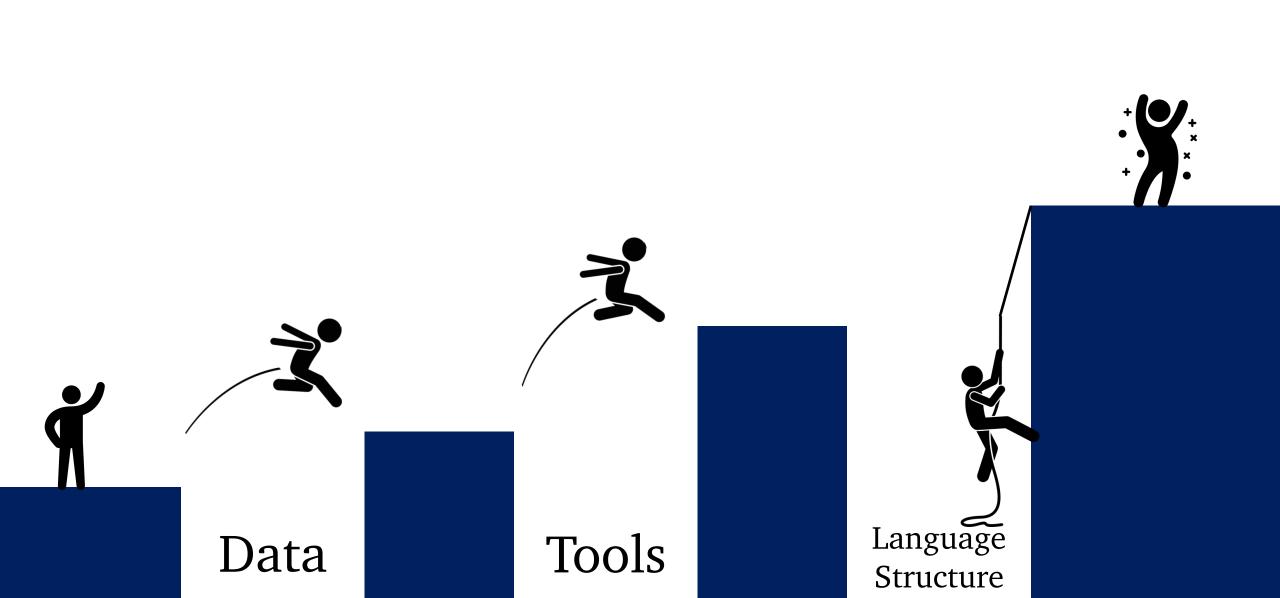
Results: Machine Translation

	Machine Translation				
	Danish→Kalaallisu <mark>t</mark>		Kalaallisut→Danish		
Method	chrF2	BLEU		chrF2	BLEU
None	44.6	3.4		61.5	15.4
BPE	56.4	8.7		<u>64.2</u>	<u>21.5</u>
Unigram	61.4	<u>10.1</u>		58.9	17.1

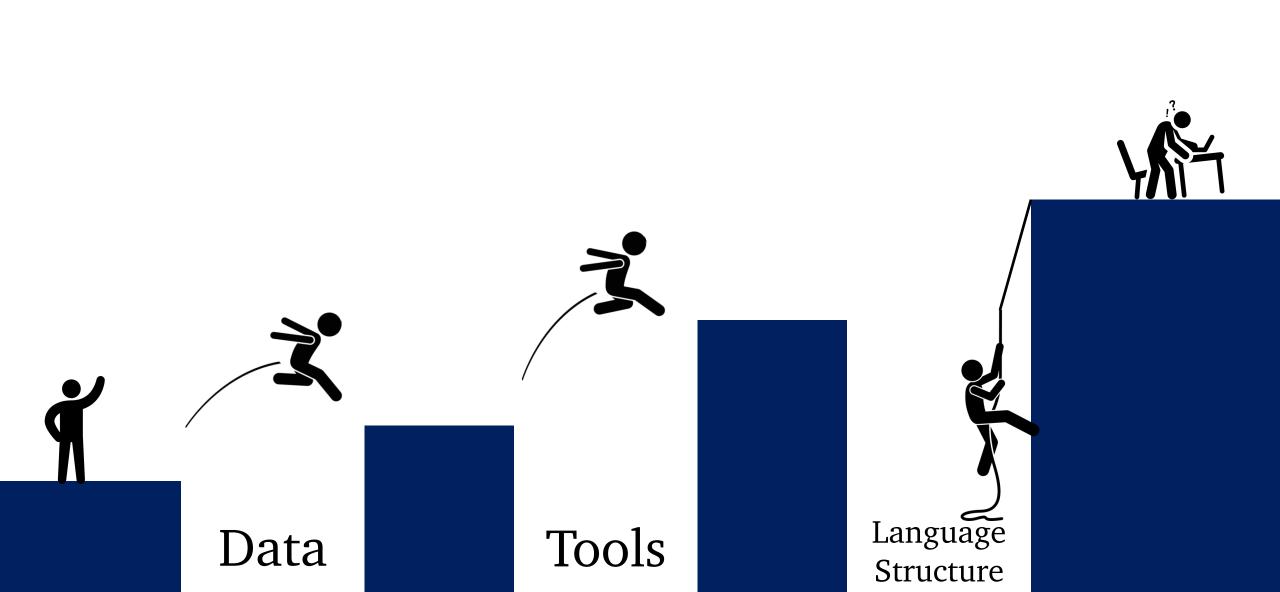
Results: Machine Translation

		Machine Translation			
	Danish-	Danish→Kalaallisut		Kalaallisut→Danish	
Method	chrF2	BLEU	chrF2	BLEU	
None	44.6	3.4	61.5	15.4	
BPE	56.4	8.7	<u>64.2</u>	<u>21.5</u>	
Unigram	61.4	<u>10.1</u>	58.9	17.1	

A Tale of Three Obstacles



A Tale of Three Obstacles



Better does not necessarily mean good...

Take-home message

The good & the bad news



• Data is (very) scarce



- Data is (very) scarce
- Out-of-the-box tools (large pre-trained language models, mainstream tokenization algorithms) are not equally applicable to all languages!



- Data is (very) scarce
- Out-of-the-box tools (large pre-trained language models, mainstream tokenization algorithms) are not equally applicable to all languages!
- Limited commercial (and academic?) interest



The Good News



The Good News

• We can be informed by structural language similarities!



The Good News

We can be informed by structural language similarities!

• Small tweaks can bring at least modest improvements!







Who will knock on a door next?

Qujanaq!

Tak!

Thank you!

And many thanks to my collaborators and the Greenlandic Language Secretariat!



Paola Saucedo



Johannes Bjerva



Ross Deans Kristensen-McLachlan



Heather Lent

Questions?



Esther Ploeger

Department of Computer Science Aalborg University, Copenhagen

espl@cs.aau.dk