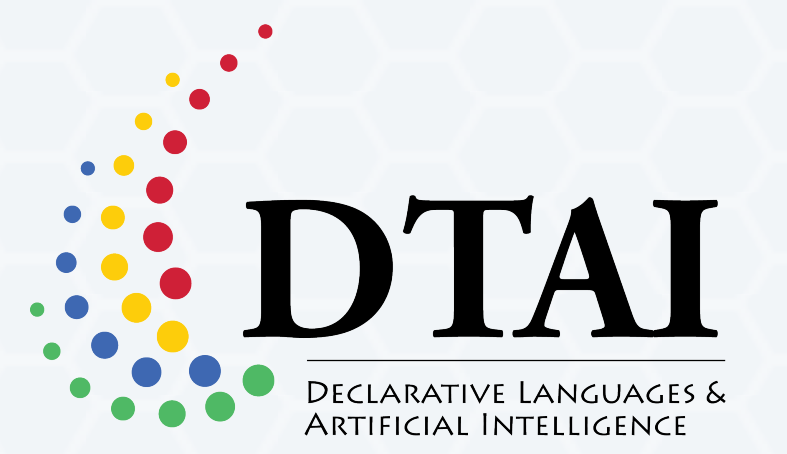


# Pruning Pre-existing BPE Tokenisers with Backwards-compatible Morphological Semi-supervision



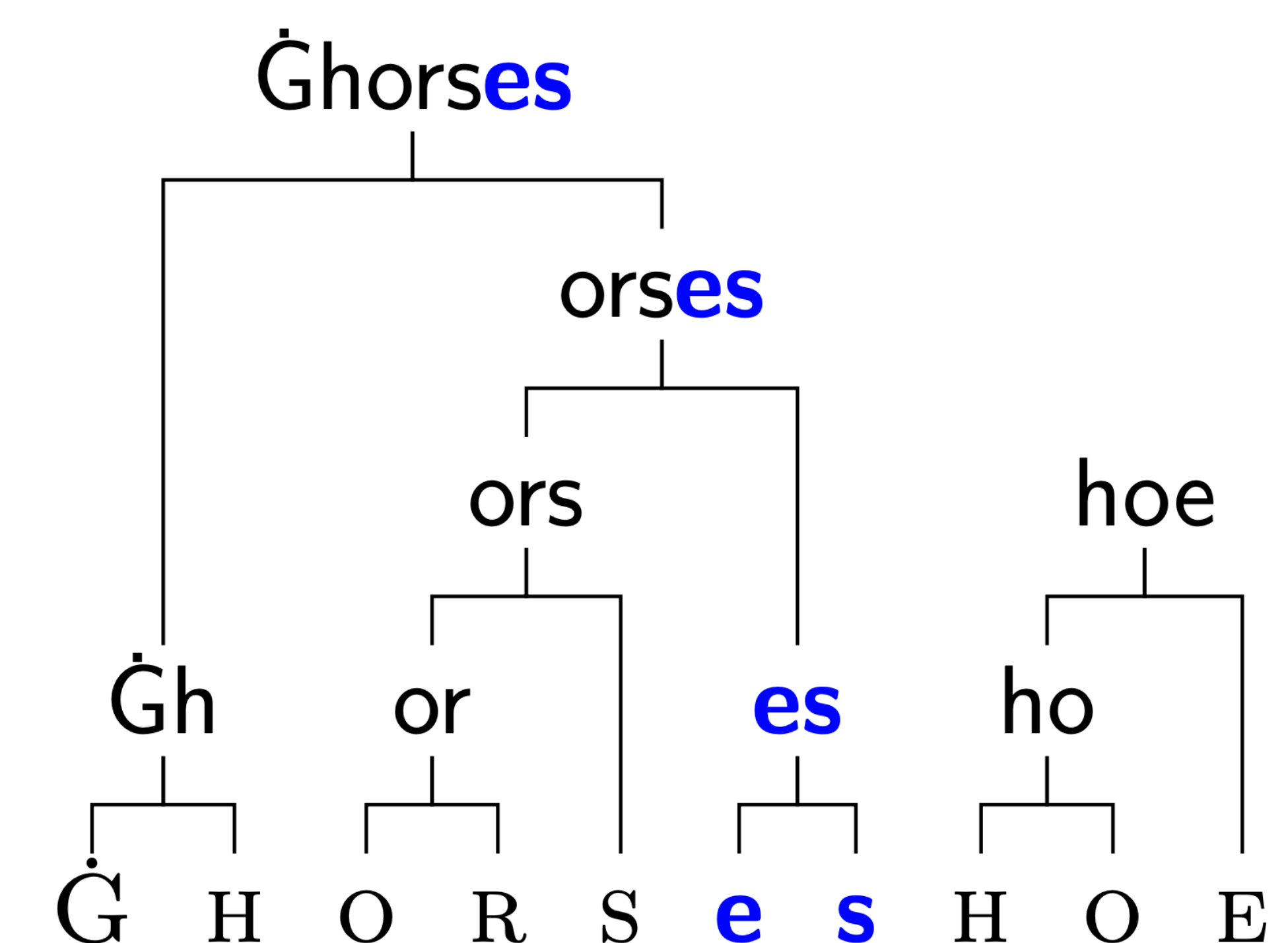
 ThomasBauwens\_

 [pieterdelobelle](#)



## Morpheme boundaries disappear

	BPE	morphology
Sennrich et al. (2016)	Gesundheits forsch  <b>ungs</b>  stitute	Gesund heit s forschungs institute
Bostrom et al. (2020)	comple t ely t ric y cles n an  <b>ote</b>  chn ology	complete ly tri cycle s nano techno logy
He et al. (2020)	vul n  <b>era</b>  bility t ighter <b>emb</b>  arked predic  <b>table</b>	vulner ability tight er embark ed predict able
Vilar et al. (2021)	bewer t  <b>ungs</b>  stru mente <b>gefan</b>  gen genommen verbrau ch  <b>spru</b>  ef standard haushalt war  <b>enab</b>  teilung not arz  <b>tau</b>  tos	bewert ung s instrumente ge fangen ge nommen verbrauch s pruef standard haushalt waren abteilung notarzt auto s

[illegible]

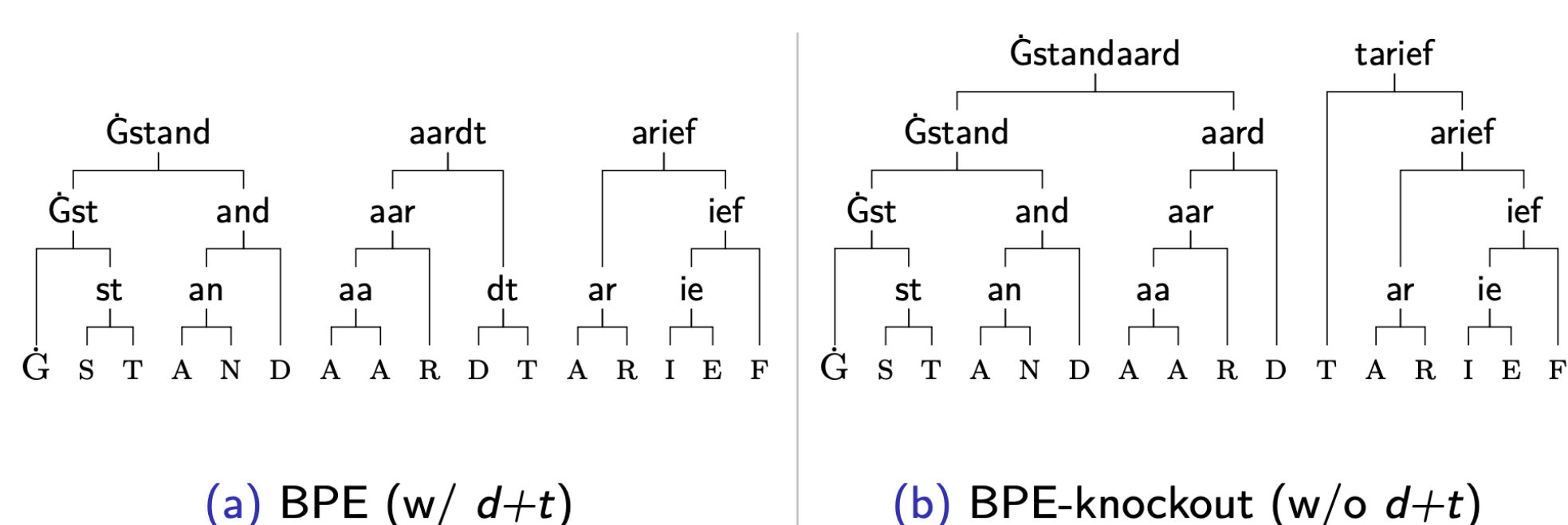
Blame merge when it contracts a morpheme boundary, then get:

$$R(m) = \frac{B(m)}{N(m)} \geq \frac{1}{2}$$

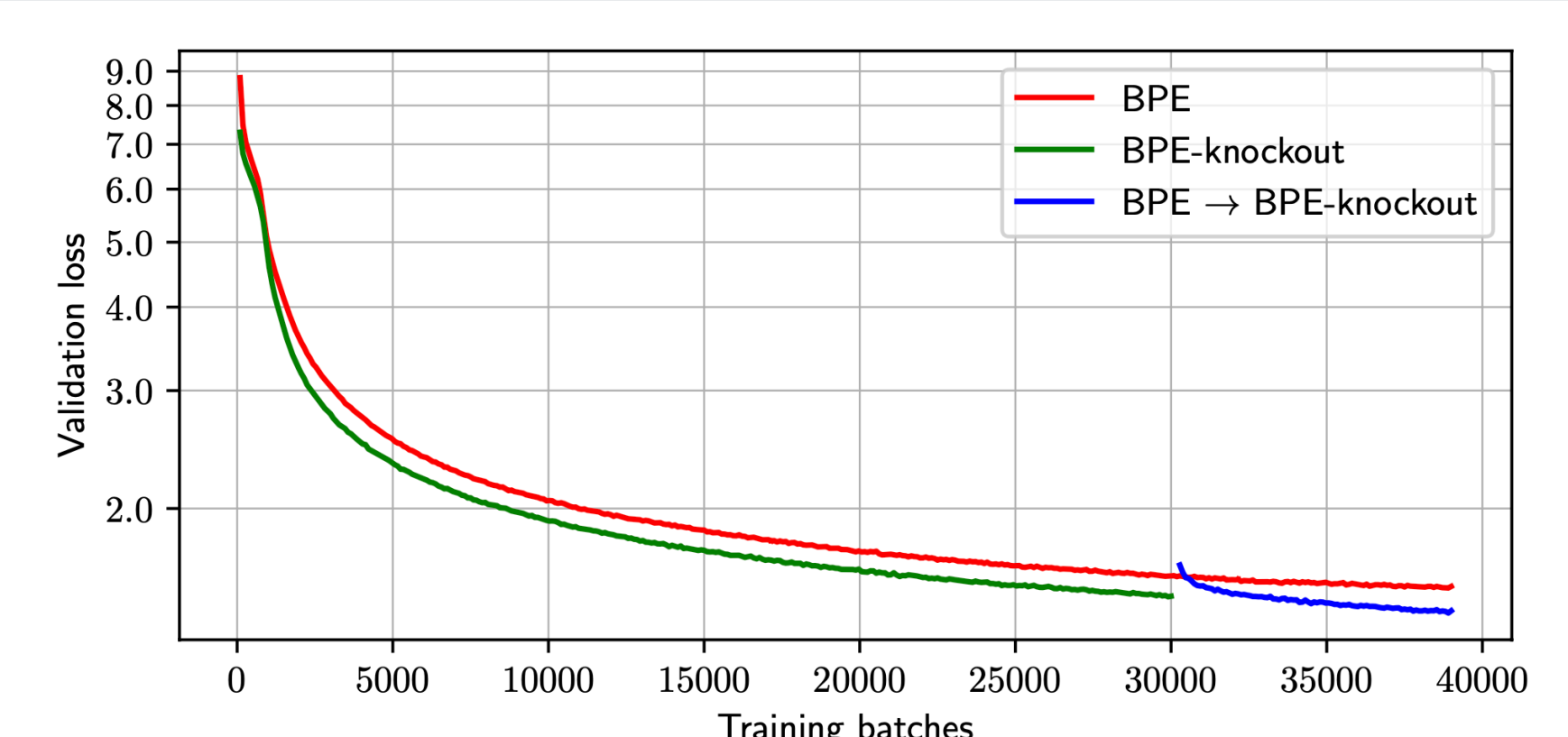
	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16					
<b>Char:</b>	d	o	c	t	o	r	a	a	t	s	m	i	s	e	r	i	e				
<b>BPE:</b>	d	o	c	t	o	r		a	a	t		s	m		i	s		e	r	i	e
<b>Gold:</b>	d	o	c	t	o	r		a	a	t		s		m	i	s	e	r	i	e	

## Switch at finetuning time = profit!

		$ V $	word types			word tokens		
			Pr	Re	$F_1$	Pr	Re	$F_1$
English	BPE	40 000	43.0	49.6	46.1	40.7	4.0	7.3
	BPE-knockout	36 814	<b>53.2</b>	<b>75.1</b>	<b>62.3</b>	<b>84.5</b>	<b>59.2</b>	<b>69.6</b>
German	BPE	40 000	45.0	54.0	49.1	54.3	8.4	14.5
	BPE-knockout	35 919	<b>55.3</b>	<b>79.8</b>	<b>65.3</b>	<b>59.5</b>	<b>69.2</b>	<b>64.0</b>
Dutch	BPE	40 000	52.6	55.3	53.9	54.3	11.0	18.4
	BPE-knockout	35 525	<b>61.7</b>	<b>78.2</b>	<b>68.9</b>	<b>81.6</b>	<b>64.1</b>	<b>71.8</b>



		PPPL		
		30k	+5k	+9k
Dutch	BPE	<b>3.88</b>	<b>3.74</b>	<b>3.66</b>
	BPE-knockout	4.67		
	BPE $\rightarrow$ BPE-knockout	200.37	4.62	4.35



		Sequence-level						Token-level					
		30k	SA +5k	+9k	30k	NLI +5k	+9k	30k	NER +5k	+9k	30k	PoS +5k	+9k
Dutch	BPE	<b>82.10</b>	81.56	81.43	<b>83.16</b>	82.55	<b>83.53</b>	80.04	83.18	83.98	93.50	85.91	93.78
	BPE-knockout	82.01			82.98			<b>86.21</b>			91.19		
	BPE → BPE-knockout	81.34	<b>82.42</b>	<b>81.74</b>	81.74	<b>82.59</b>	83.14	80.58	<b>86.77</b>	<b>87.51</b>	<b>94.82</b>	<b>96.03</b>	<b>96.01</b>

A BPE tokenizer can be improved after a model has been pre-trained on it, but more research is needed to exploit morphology in tokenizers.



Learn more at  
*pieter.ai/*  
***bpe-knockout***