

Mixed Orthographic/Phonemic Language Modeling: Beyond Orthographically Restricted Transformers (BORT)

Robert C. Gale
galer@ohsu.edu

Alexandra C. Salem
salem@ohsu.edu

Gerasimos Fergadiotis
gf3@pdx.edu

Steven Bedrick
bedricks@ohsu.edu

NLP for SLP

In the field of speech language pathology, professionals work to diagnose and treat language disorders. Transcripts from people with language disorders often include speech errors transcribed with phonemes. In order to automate assessment we need models that can represent phonemes.

BORT

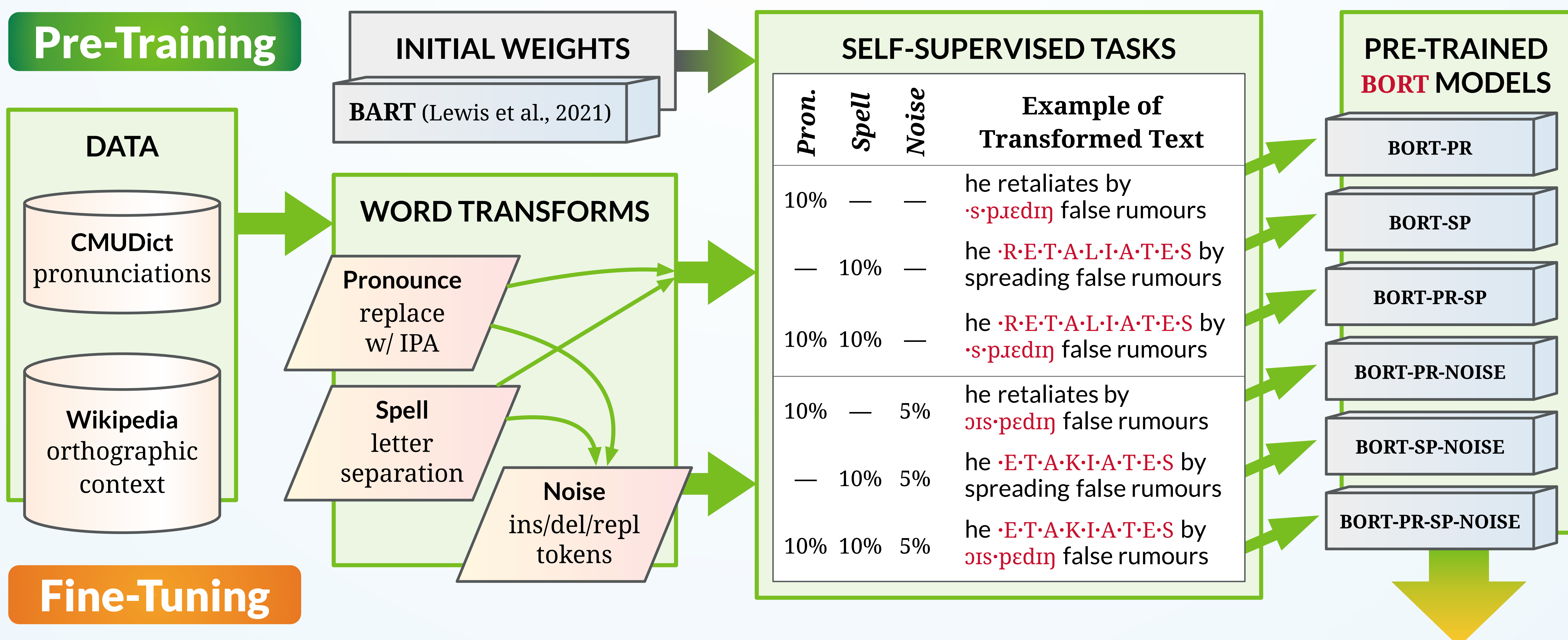
is a pre-trained LLM for **clinical language assessment** designed for use on mixed **phonemes & orthography**

This work was funded by NIH/NIDCD award #R01DC015999 awarded to Steven Bedrick and Gerasimos Fergadiotis

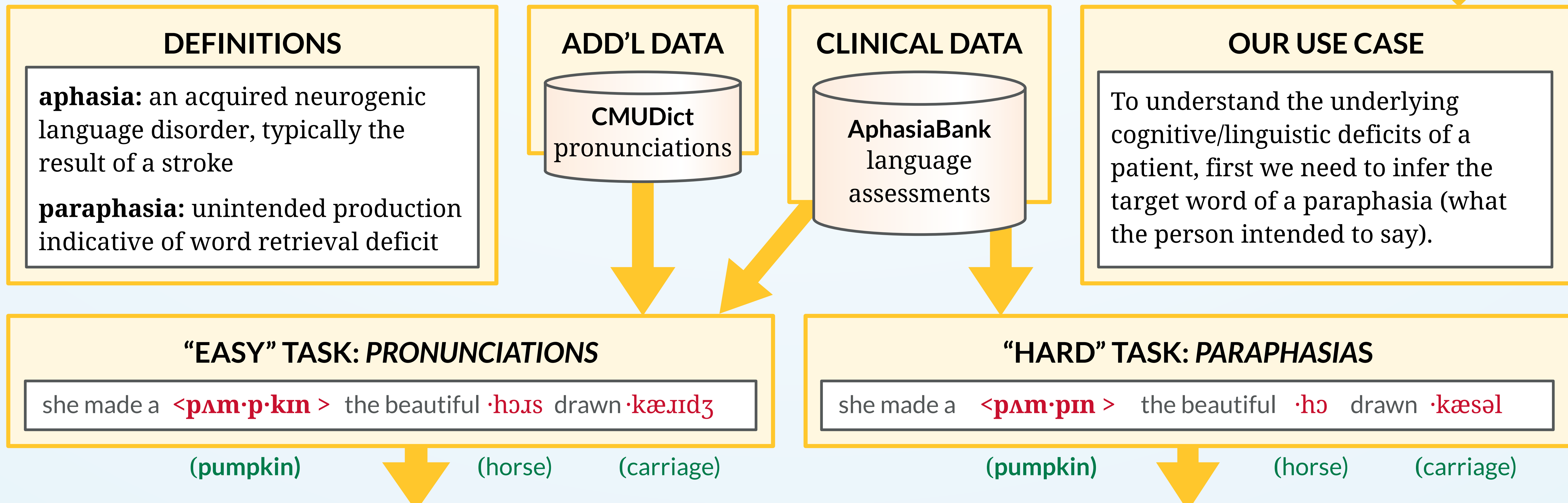


Acknowledgements

Pre-Training



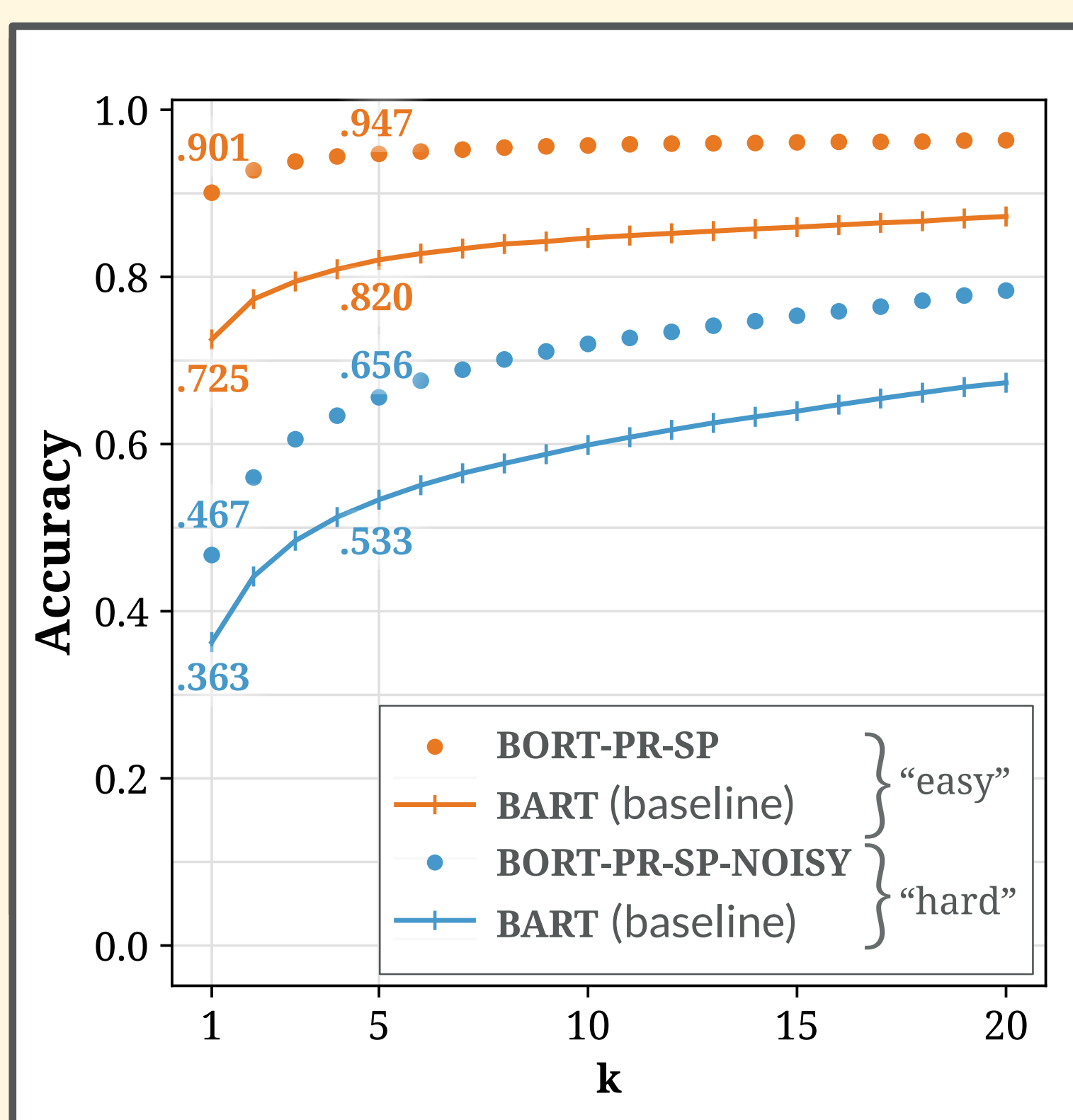
Fine-Tuning



"EASY" RESULTS

Pre-training Configuration	CER		Accuracy	
	Top 1	Top 1	Top 1	Top 5
BORT-PR	.083	.869	.931	
BORT-SP	.106	.843	.906	
BORT-PR-SP	.057	.901	.947	
BORT-PR-NOISE	.089	.863	.925	
BORT-SP-NOISE	.096	.848	.911	
BORT-PR-SP-NOISE	.060	.895	.947	
BART-BASE	.228	.725	.820	

TOP-K ACCURACY



"HARD" RESULTS

Pre-training Configuration	CER		Accuracy	
	Top 1	Top 1	Top 1	Top 5
BORT-PR	.462	.451	.634	
BORT-SP	.526	.401	.579	
BORT-PR-SP	.447	.456	.641	
BORT-PR-NOISE	.452	.458	.640	
BORT-SP-NOISE	.469	.446	.625	
BORT-PR-SP-NOISE	.420	.467	.656	
BART-BASE	.606	.363	.533	

CER = character error rate

Boldface values were significantly different from all other models, with the exception of those italicized, according to McNemar's test.

MODELS & RESOURCES:

<https://github.com/rcgale/bort>

