Multimodality extension to Universal Multilingual BPE Text Tokenizer

This abstract is proposing Multimodality extension to the paper One Tokenizer To Rule Them All..[1]. The referenced paper [1] mainly uses bucketed weighting scheme on unseen/expanded languages by script e.g. (Devnagari, Hindi) or (Latin, Polish) pair and trains Byte-Pair Encoding (BPE) model on diverse text corpus across \sim 69 languages (combination of languages used in pretraining and many others that are only intended for tokenizer coverage). It also provides byte-fallback for edge cases outside training data. This abstract is about concrete enhancements to the above paper's [1] bucketed weighting scheme and explicitly accounts for multiple modalities like Images, Speech, OCR, Text etc. The below enhancements are aimed at achieving >= 0.95 CMS score to consider the resultant tokenizer as Multimodal tokenizer. 1.Modality-aware buckets[2] Extend buckets to [script, modality] pairs e.g., (Devanagari, OCR), (Arabic, ASR) and assign higher sampling weights to underrepresented pairs and monitor per-bucket coverage in tokens/word and bytes/token. Keeping total vocabulary same; train BPE on weighted samples. Measure improvement in OCR/ASR-related tasks for same scripts; per-bucket tokens/word. Success: ≥ small positive lift (0.5–2 pts) on OCR-heavy tasks for those scripts vs baseline. 2. Confidence-weighted sampling Use OCR/ASR confidence scores to downweight low-confidence examples or to preferentially sample medium-confidence ones for tokenizer training. Integrate OCR/ASR confidence scores and sample with prob \propto (α + conf^{6}). Precompute confidences; tune α (e.g., 0.05) and β (e.g., 1.0 \rightarrow 2.0). Keep some low-confidence included via α. Measure: Token noise (tokens seen only in low-confidence data), downstream VQA/DocVQA on OCR; stability of merges. Success: Cleaner merges (fewer spurious tokens) and small downstream improvement; reduction in tokens primarily seen in noisy buckets. 3. Adaptive Reweighting with Feedback [3] During tokenizer training, periodically evaluate downstream proxy tasks (small VQA/ASR validation slices). Reweight buckets that show poor downstream performance. Every N steps, compute per-bucket validation loss; increase sampling weight for buckets with high loss (up to a cap). Measure: Convergence speed on proxies; stability of vocab. Success: Faster improvements on held-out proxies; stable vocabulary.4. Cross-Modal Coverage Balancing[4] Upweight text segments that are aligned to images/speech (e.g., OCR region + image) so merges capture visuallygrounded tokens by marking multimodal aligned text and multiply sample weight by γ (1.5–3.0) and measure improvement in grounded retrieval/DocVQA EM and fewer mis-OCR tokens. Success: Noticeable lift on grounding tasks (>= 1-3 pts) 5. Curriculum-Based Bucket Scheduling[5] Systematic phases in the training as in Phase A, clean high-confidence multimodal pairs. Phase B: gradually add noisy/augmented examples for N steps. Measure: Merge stability (fewer reversions), downstream robustness to noisy OCR. Success: Better OCR robustness and fewer lowquality tokens. 6. Multimodal-Aware Validation Metrics[6][7][8][9] We want metrics that reflect multimodal performance, not just perplexity or compression. Composite Multimodal Score (CMS) should be >= 0.95 AND no single Primary Metric falls below target.

Metric (Mi)	Dataset(s)	Expected Value /	Priority
Medic (Mi)	Da(ase((s)	Success Target	Filolity
Text Perplexity	WikiText, Flores, mC4	Within 5% of baseline	Primary
restreiplesity	(clean text)	tokenizer	rilliary
A 7.			D :
Avg. Tokens per	WikiText, Flores	≤ 1.1x baseline	Primary
Word (Text CR)			
OCR	FUNSD, SROIE, IIT-	≤ 1.3x text CR	Primary
Compactness	CDIP, SynthText	l	I I
Ratio (OCR-CR)			
OCR Robustness	FUNSD, SROIE	≥ 90% overlap with	Secondary
Score (Edit		GT tokens	I I
distance overlap)			
ASR Token Error	LibriSpeech	Should≋ WER (gap≤	Primary
Rate (TER)	(clean/other),	3%)	
WER-to-TER	LibriSpeech, MLS	≤+3% gap	Secondary
Gap	·		l 'I
DocVQA Exact	DocVQA, TextVQA	≥ 95% of baseline EM	Primary
Match (EM)		l	'
DocVQAF1	DocVQA	≥ 95% of baseline F1	Secondary
Score		l	l 'I
Caption	COCO Captions,	Within 10% of	Primary
Compactness	Flickr30k	baseline tokenizer	'
(tokens per		l	I I
caption)		l	l
Alignment	COCO Captions + CLIP	≥ baseline tokenizer	Secondary
Overlap	labels	IOU	I 'l
(IOU@token)		l	I I
Composite	All above datasets	Overall≥ 0.95 relative	Primary
Multimodal Score		to baseline	1 ~ 1
(weighted across			I I
modalities)		l	1
		•	

Mi = measured metric value (Refer Metric column from side table) Ti = target threshold (so ratio ≥1 is "good") wi = weight assigned to modality (wi values = [Core text = 20%, OCR = 25%, ASR = 25%, DocVQA = 15%, Captions = 15%]) References -[1] https://arxiv.org/pdf/2506.10766 [2] https://www.rohanpaul.com/p/multilingual-and-multimodal-llms [3] https://aclanthology.org/D17-1158/ [4] https://arxiv.org/abs/1909.11740 [5] https://dl.acm.org/doi/10.1145/1553374. 1553380[6] https://arxiv.org/abs/2007.00398 [7]https://arxiv.org/abs/1904.08920[8] https://www.openslr.org/12[9] https://commonvoice.mozilla.org/en Tools used – ChatGPT for validating the approaches and finding the

CMS = ∑ wi · Mi/Ti

Tools used – ChatGPT for validating the approaches and finding the references to the proposed approaches in this abstract for citation purpose.