

Multimodal Pretraining Unmasked: A Meta-Analysis and a Unified Framework of Vision-and-Language BERTs

Emanuele Bugliarello Ryan Cotterell Naoaki Okazaki Desmond Elliott

TACL / EMNLP 2021

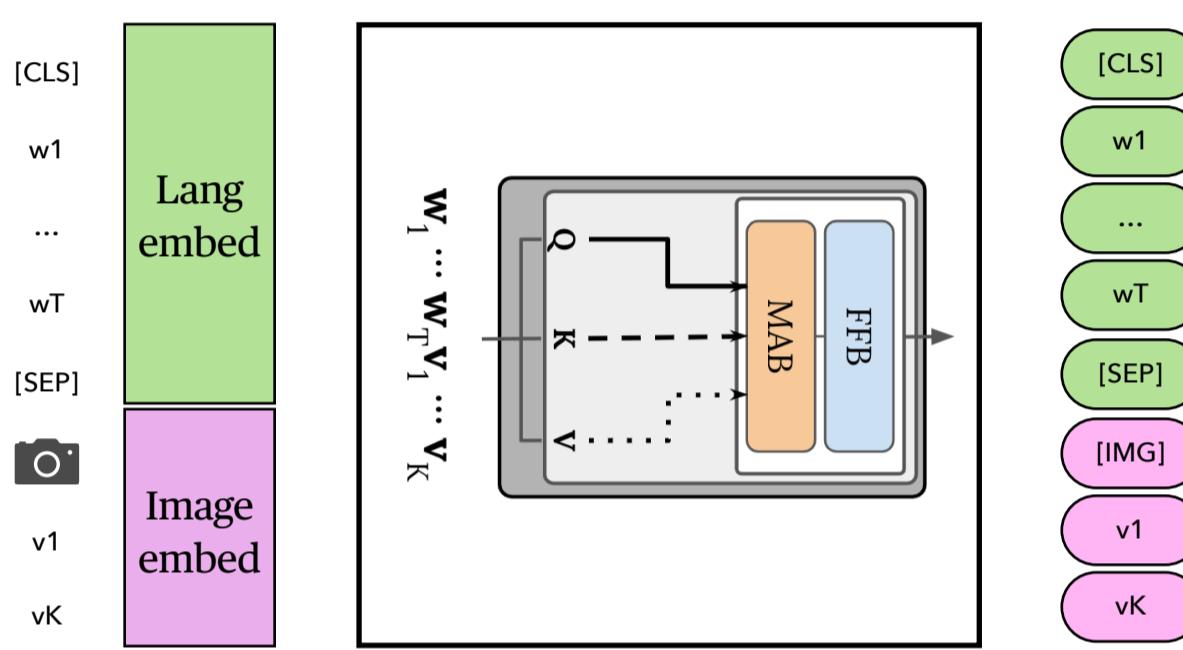
UNIVERSITY OF
OPENHAGEN



V&L BERT Families

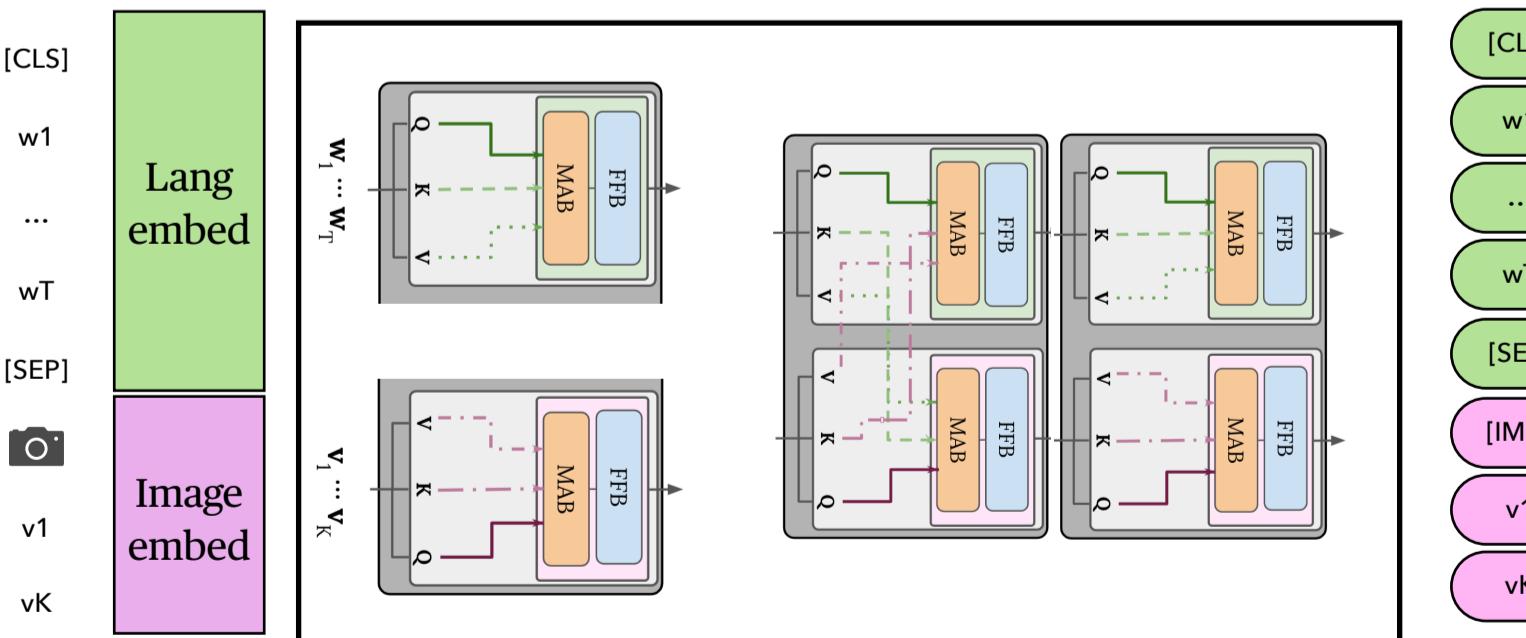
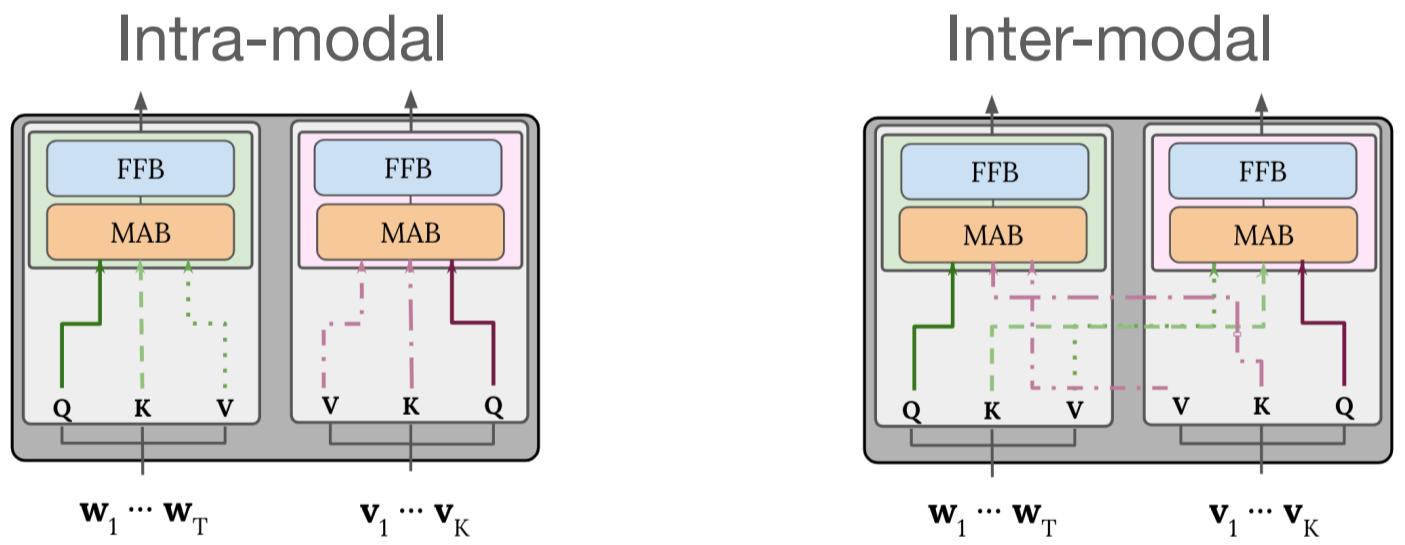
Single-Stream

- Concatenate image—text inputs



Dual-Stream

- Model the Image and text independently
- Two types of cross-modal layers

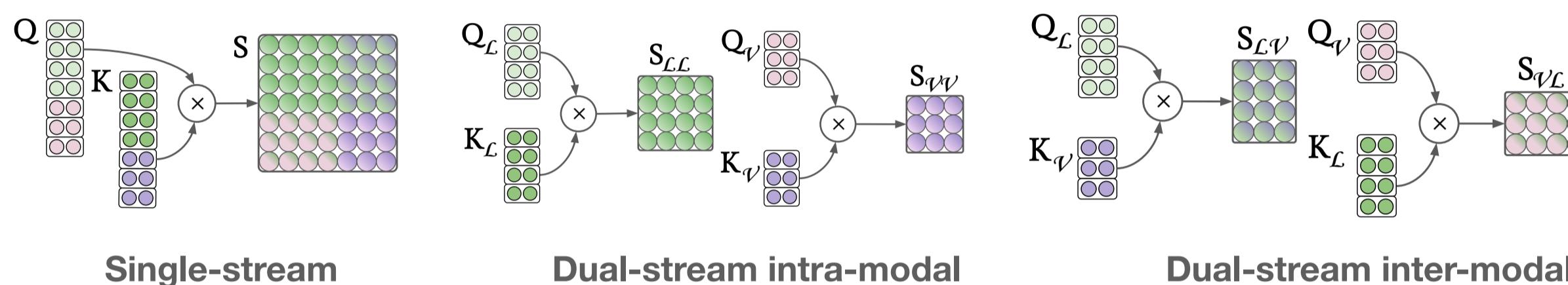


Code, Data & 50 Pretrained Models

- github.com/e-bug/volta
- github.com/e-bug/mpre-unmasked

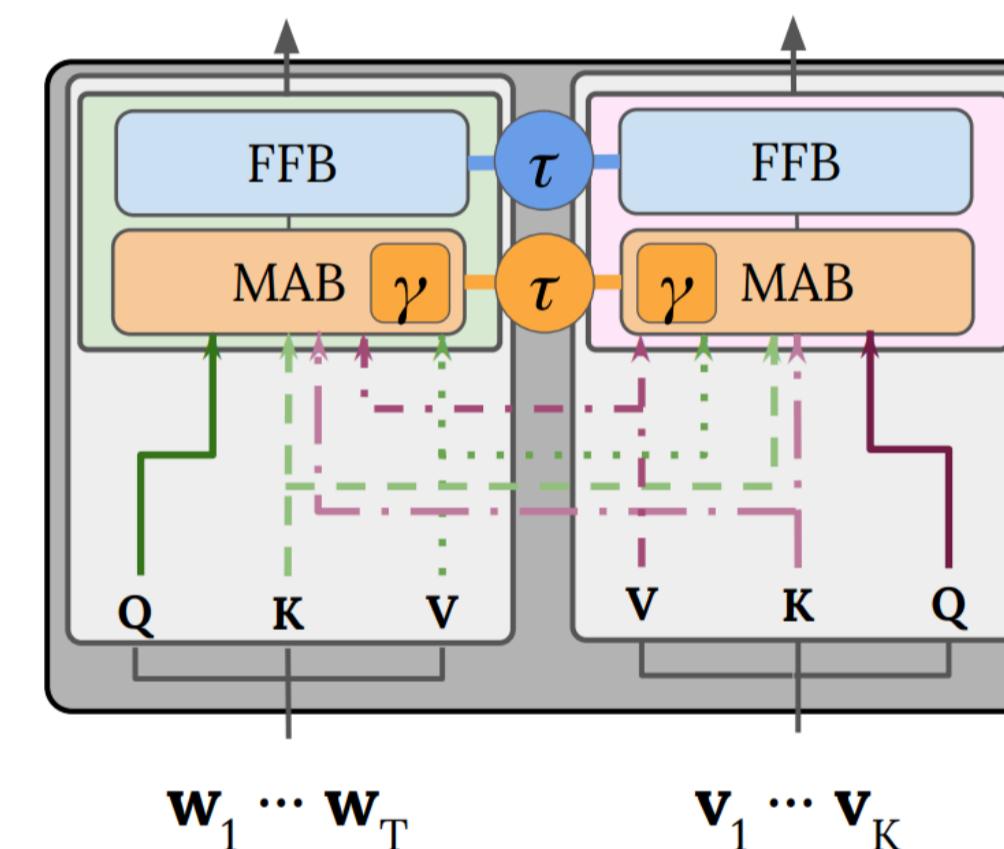
VOLTA: A Unified Framework

Attention matrices

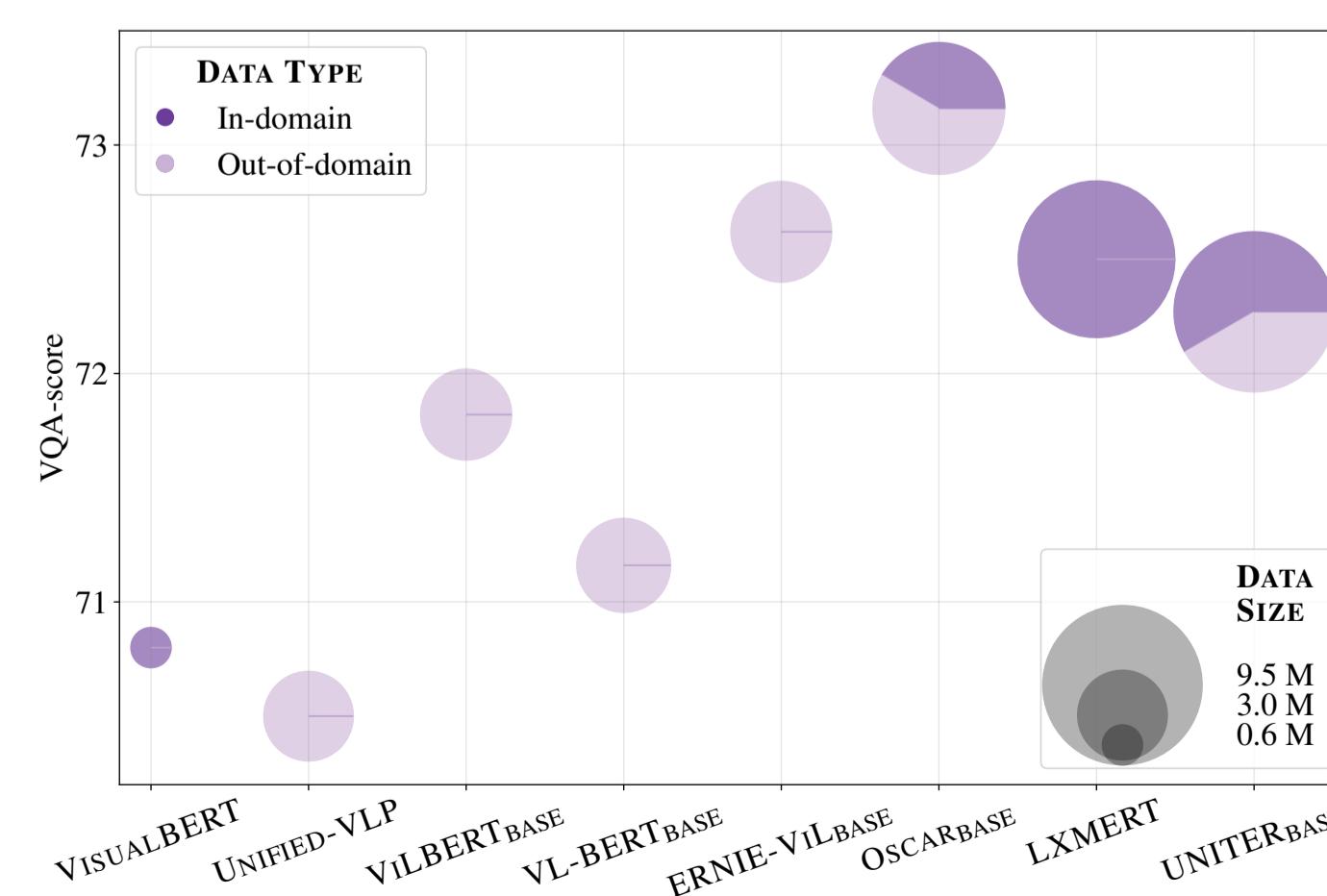


New Gated Bimodal Transformer Layer

- Express many architectures with this layer**
- Single- & dual-stream layers are special cases
- Takes a set of fixed binary variables $\{\gamma, \tau\}$
 - γ : regulate the cross-modal interactions
 - τ : control if parameters are tied between modalities



Which V&L BERT? 🤔



These models are trained on different amounts & types of data

- In-domain:** also relevant for downstream tasks (e.g. COCO)
- Out-of-domain:** unrelated to downstream tasks (e.g. CC)

Controlled Experiments

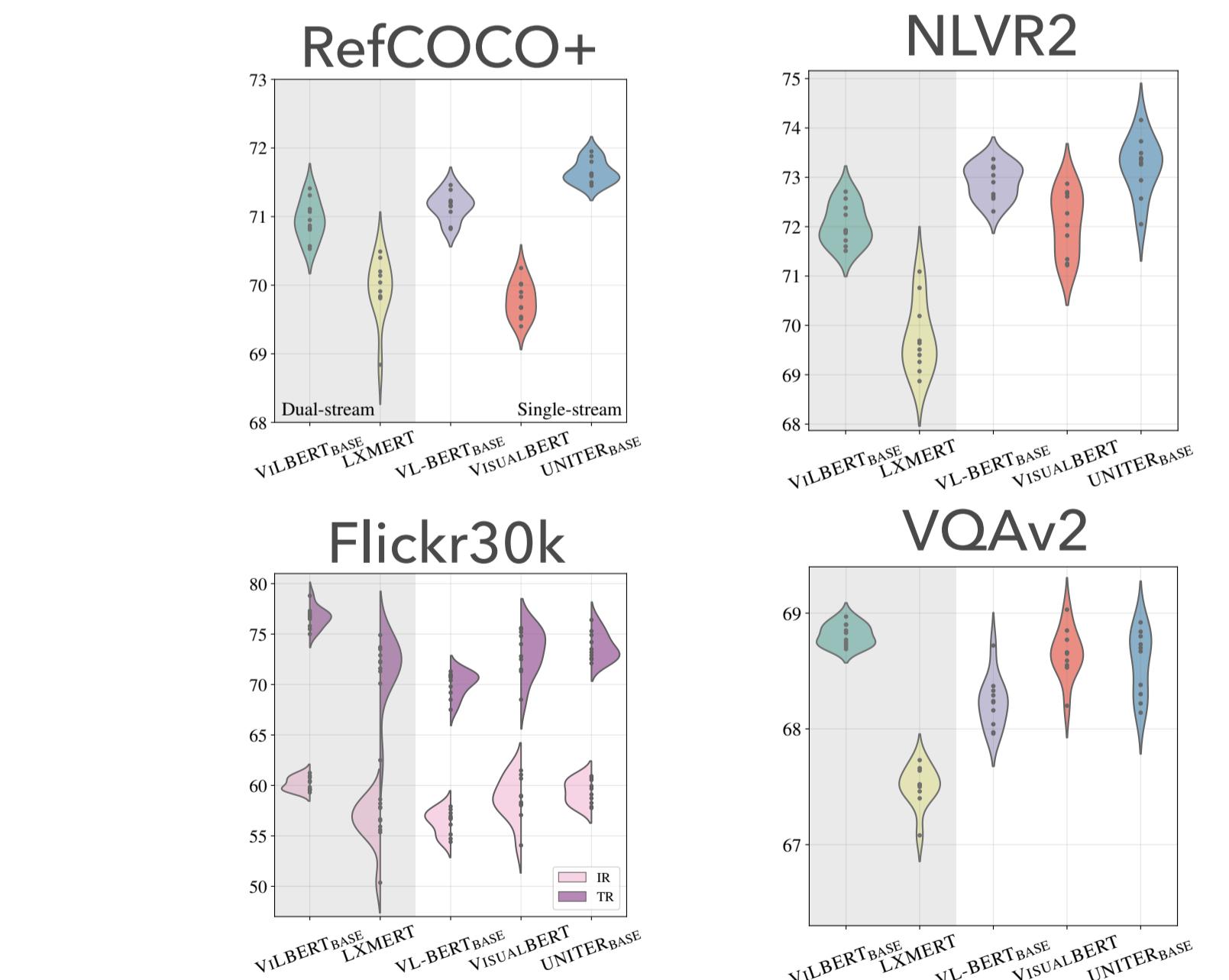
- We implement existing models in VOLTA
- Train them on the same setup with multiple seeds and the MLM + MRC-KL + ITM objectives

Models:

- Single-stream: VL-BERT, VisualBERT, UNITER
- Dual-stream: ViLBERT, LXMERT

Pretraining data: Conceptual Captions (2.77M)

Results



- Substantial variation due to seed
- Similar performance** in controlled setup
- Single- & Dual-Stream are **on par**
- Embedding layer** is crucial