

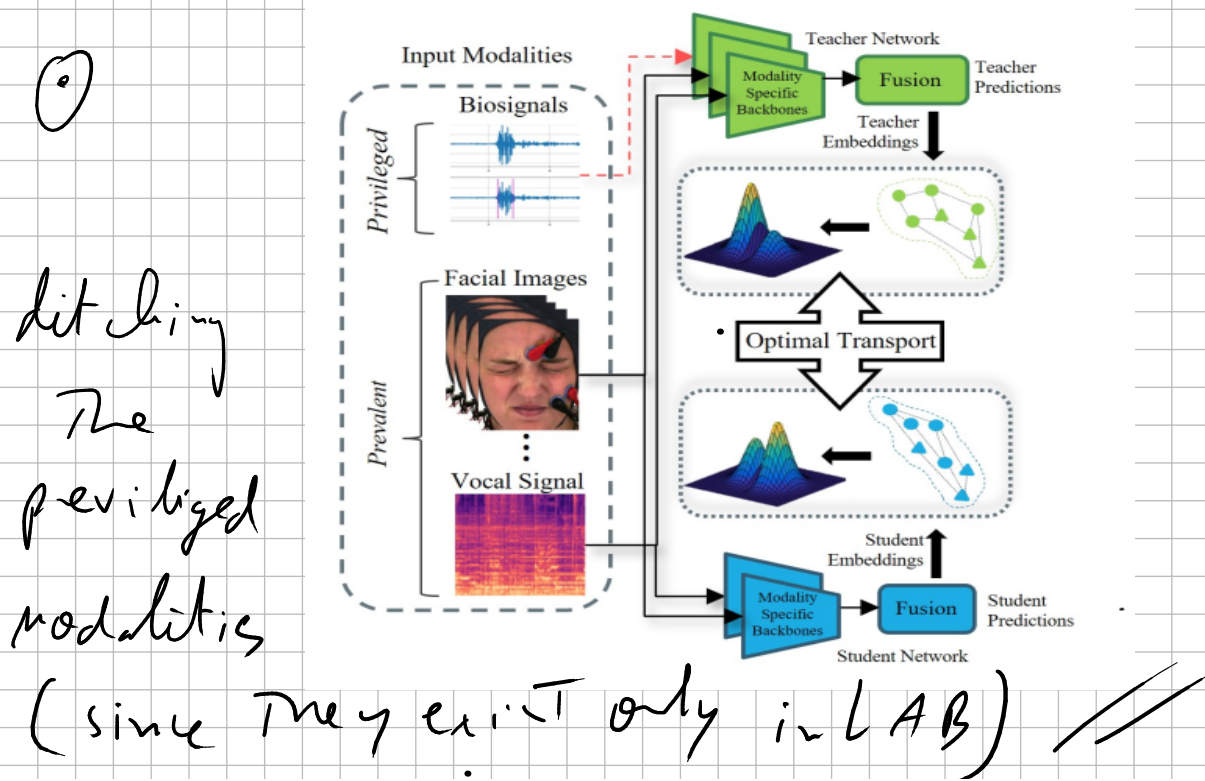
* KD optimal Transport Paper *

① Multi-modal expression recognition is a problem solved long ago, in Labs; once these models are put in the wild, they struggle since not all modalities are present at Test Time (vision, speech, bio-signals), in practice, bio-signals are pretty much absent. (This type of info is called privileged)

⇒ KD have been proposed as a way to distill info from multiple Teacher models, each trained on a modality, to one student model. But, KD uses point to point matching ⇒ student has to match each modality separately ⇒ has no way of capturing cross-modality structure.

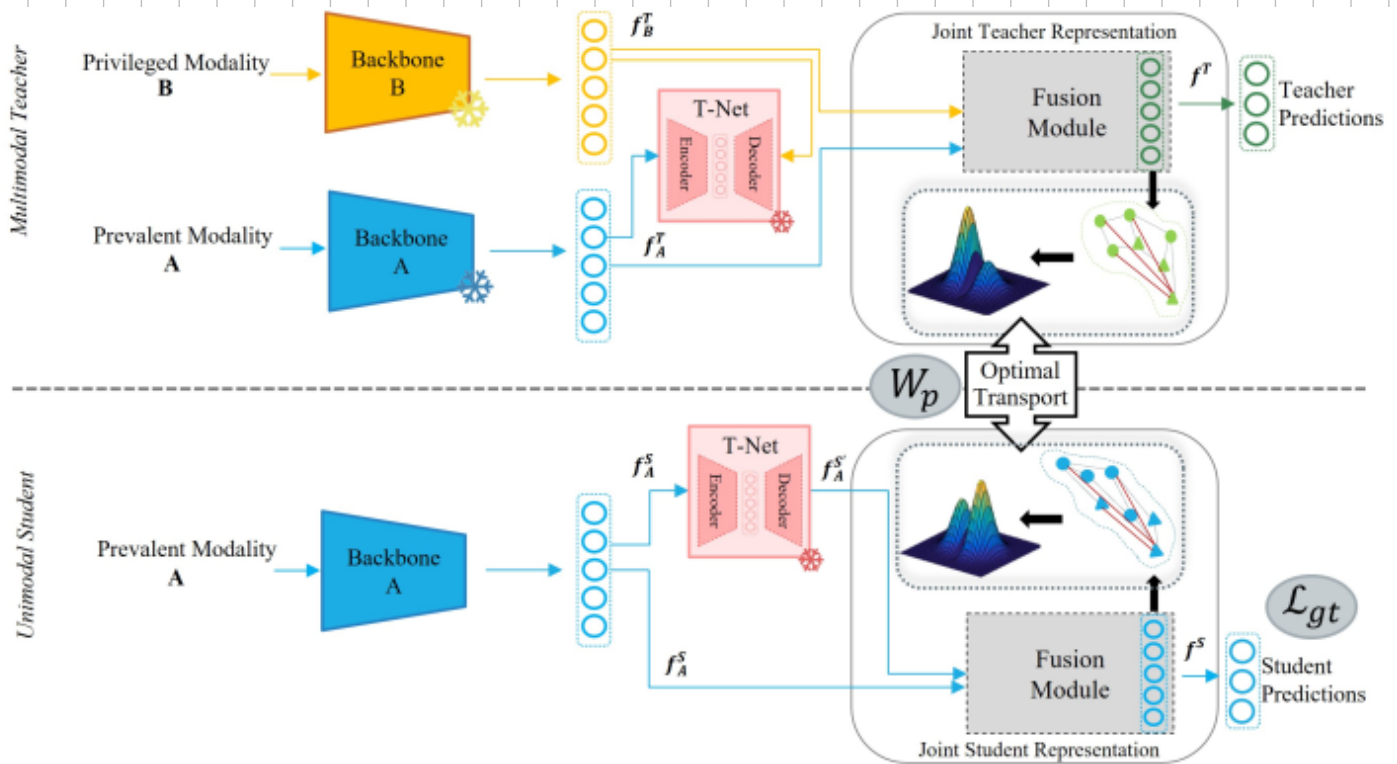
⇒ This paper proposes a way for The student Model to capture that structure & argues that it boosts student performance

① KDOT capture the cross-modality structure by calculating cosine similarity & selecting the top- k anchors (to allow for a flexible distillation process, not rigid)



① PKNOT's goal is To mitigate The shortcomings privileged Modalities introduce To The model (dipped performance at inference Time)

② architecture of PKNOT :



Teacher Network is Trained on both prevalent & privileged data To boost performance, & it's Trained normally after, (Mid-fusion Paradigm);

a very interesting block is introduced;
That's The "T-Net", whose role is To map (prevalent, privileged)
modality for each sample, it is trained with The
Teacher Network & is used "inference Mode" when
The distillation is happening, it's role is To hallucinate
privileged Modalities based on prevalent ones; a good way To
create dummy privileged data That's gonna be used To
reinforce The cross-modality structure That's gonna be
distilled from Teacher To Student.

- Teacher-Student distillation uses OT Method & selects
k-anchors (most dissimilar samples in The batch)
To distill based on. (will be explained thoroughly later).

⊗ Structural similarity matching (KNOT + k-anchors)

