



Carmelo Scribano, Davide Sapienza, Giorgia Franchini, Micaela Verucchi, Marko Bertogna Department of Physics, Informatics and Mathematics, University of Modena and Reggio Emilia, Modena (Italy)

> The AI City Challenge Track 5: Contribution

All You Can Embed (AYCE), a modular solution to correlate single-vehicle tracking sequences with natural language. The main building blocks of the proposed architecture are (i) BERT to provide an embedding of the textual descriptions, (ii) a convolutional backbone along with a Transformer model to embed the visual information, (iii) a new loss function designed for the track.

> Natutal Language Branch (BERT Finetuning)





All You Can Embed: Natural Language based Vehicle Retrieval with Spatio-Temporal Transformers



\succ The objective function

 $\mathcal{L}_{AYCE}(A, P, N) = TL(A, P, N) +$

Where:

$$\Phi(A^i, P^i) = \min\left(d\left(A^i_{m,}P^i_n\right)\right) \ m, n \in \{1, 2, 3\}$$

TL is defined as the standard Triplet Loss, the distance measure used in TL is the average of all distance pairs (3)(3)Language Visual between and outputs.

Architecture Overview





$$+\frac{1}{Bs}\sum_{i=1}^{Bs}\beta\cdot\Phi(A^{i},P^{i})$$





Team	MRR
Alibaba	0.1869
TimeLab	0.1613
SBUK	0.1594
UNIMORE	0.1078

https://github.com/cscribano/AYCE_2021