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

AI City Challenge 2021

Carmelo Scribano[1,2], Davide Sapienza[1,2], Giorgia Franchini[1,3]
Micaela Verucchi[1] and Marko Bertogna[1]

[1]University of Modena and Reggio Emilia, [2]University of Parma, [3]University of Ferrara

# Problem statement

"*Natural language (NL) description offers another useful way to specify vehicle track queries. In this new Challenge Track, participating teams will perform vehicle retrieval given single-camera tracks and corresponding NL descriptions of the targets.*"

*- AI City Challenge 2021 – Track5*

**GOAL:**

- Matching Single-Vehicle Tracking Sequences with

  the corresponding Natural Language Descriptions.

**Related tasks:**

- Image and Video retrieval

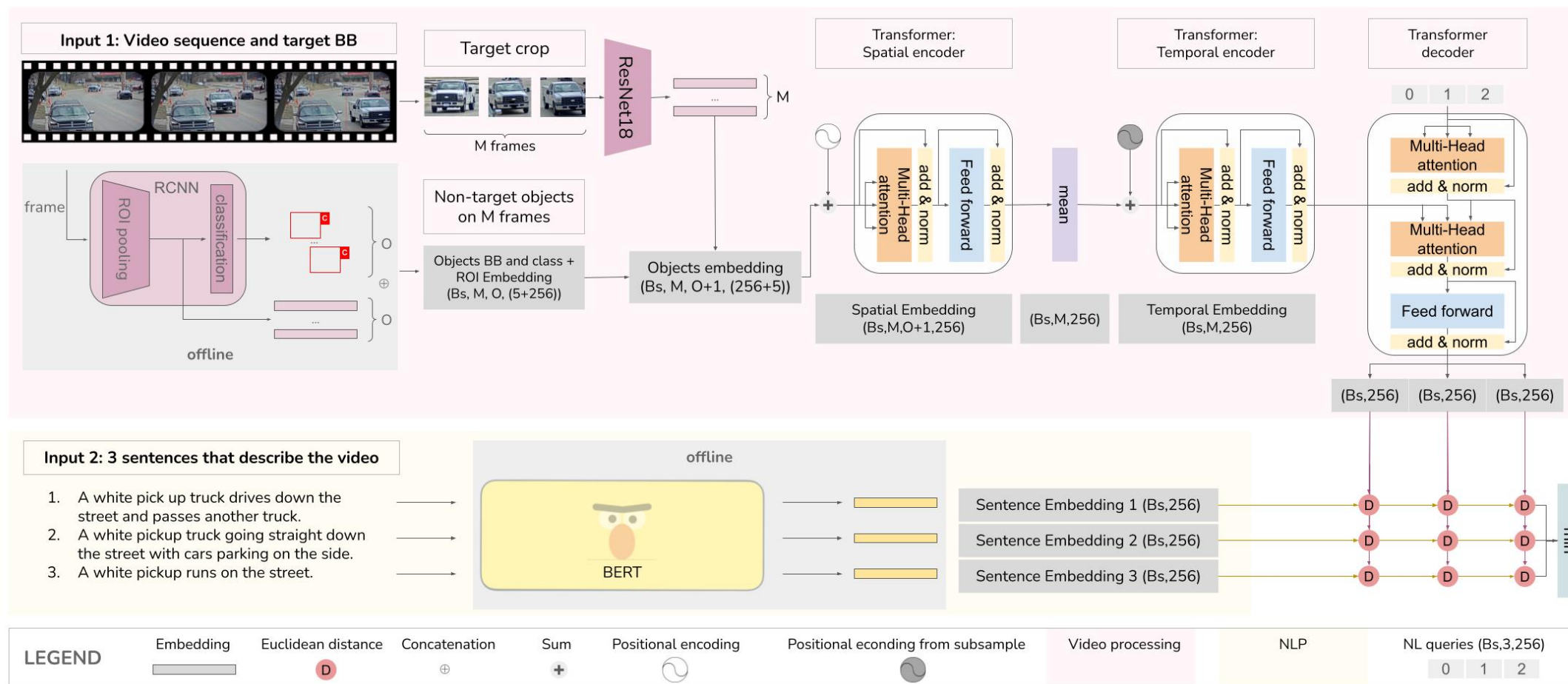- Multi-Modal video Understanding



1. A white pick up truck drives down the street and passes another truck.
2. A white pickup truck going straight down the street with cars parking on the side.
3. A white pickup runs on the street.

1. A red SUV runs down the street alongside parked cars.
2. Red SUV keeps straight followed by a maroon car.
3. A red SUV runs down the road followed by a black vehicle.
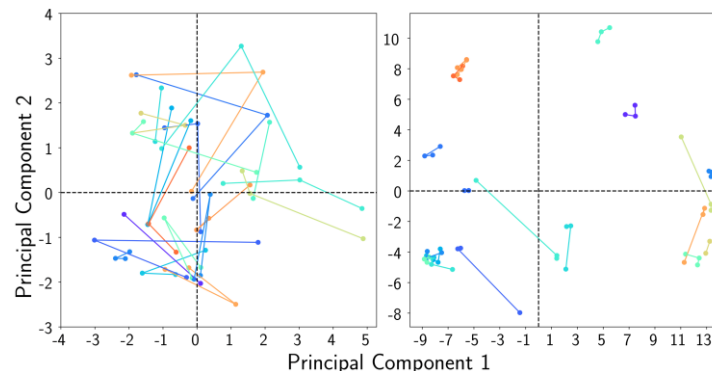
# Model Overview

# Natural Language Branch

**OBJECTIVE:**

$$\begin{cases} d(T_1^i, T_2^i) \ll d(T_1^i, T_2^j) \\ d(T_1^i, T_2^i) \approx 0 \end{cases} \quad \forall i : i \neq j$$

**Where:**

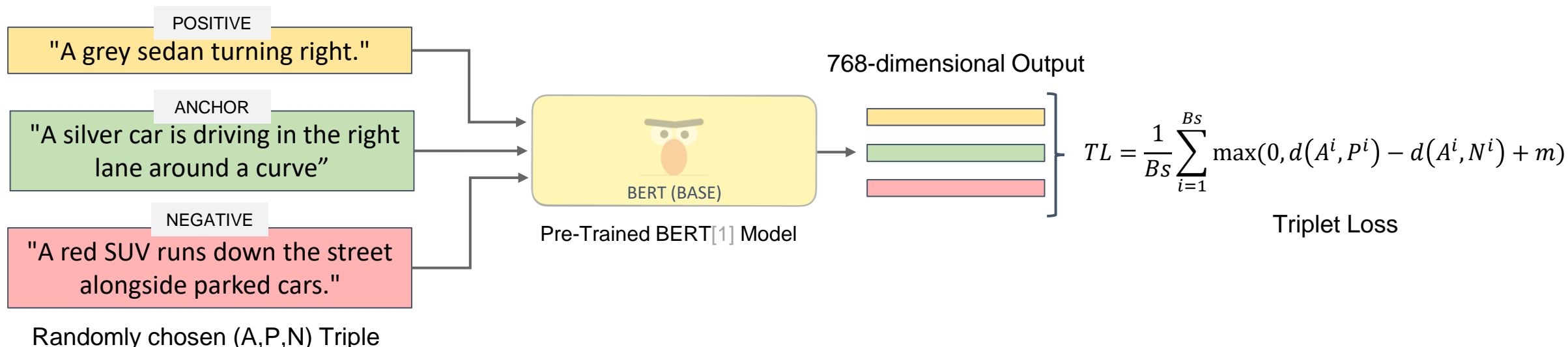$$T_a^i = BERT(t_a^i) \qquad d(u,v) = 1 + \frac{u^T v}{||u|| ||v||} \in [0,2]$$

**BEFORE:**

$$mean(d_{INTER}) = 0.1703$$
$$mean(d_{INTRA}) = 0.1899$$

**AFTER:**

$$mean(d_{INTER}) = \mathbf{0.2089}$$
$$mean(d_{INTRA}) = \mathbf{0.6140}$$

POSITIVE

"A grey sedan turning right."

ANCHOR

"A silver car is driving in the right lane around a curve"

NEGATIVE

"A red SUV runs down the street alongside parked cars."

Randomly chosen (A,P,N) Triple

768-dimensional Output

BERT (BASE)

Pre-Trained BERT[1] Model

$$TL = \frac{1}{Bs}\sum_{i=1}^{Bs} \max\left(0, d(A^i, P^i) - d(A^i, N^i) + m\right)$$

Triplet Loss

# Visual Branch



Object
Embeddings

Original Sequence
(1..3620)

Uniform Sampling
**N < 80**

$(V_1, V_2, V_3)$

Evaluation Mode
(Inference)

$$d(V, T) = \min\Big(d\big(V_i, T_j\big)\Big)$$
$$\forall\, (i, j)$$

$t_1$
$t_2$
$t_3$

BERT (BASE)

$(T_1, T_2, T_3)$

Language Branch
*After Finetuning

# Object Embeddings

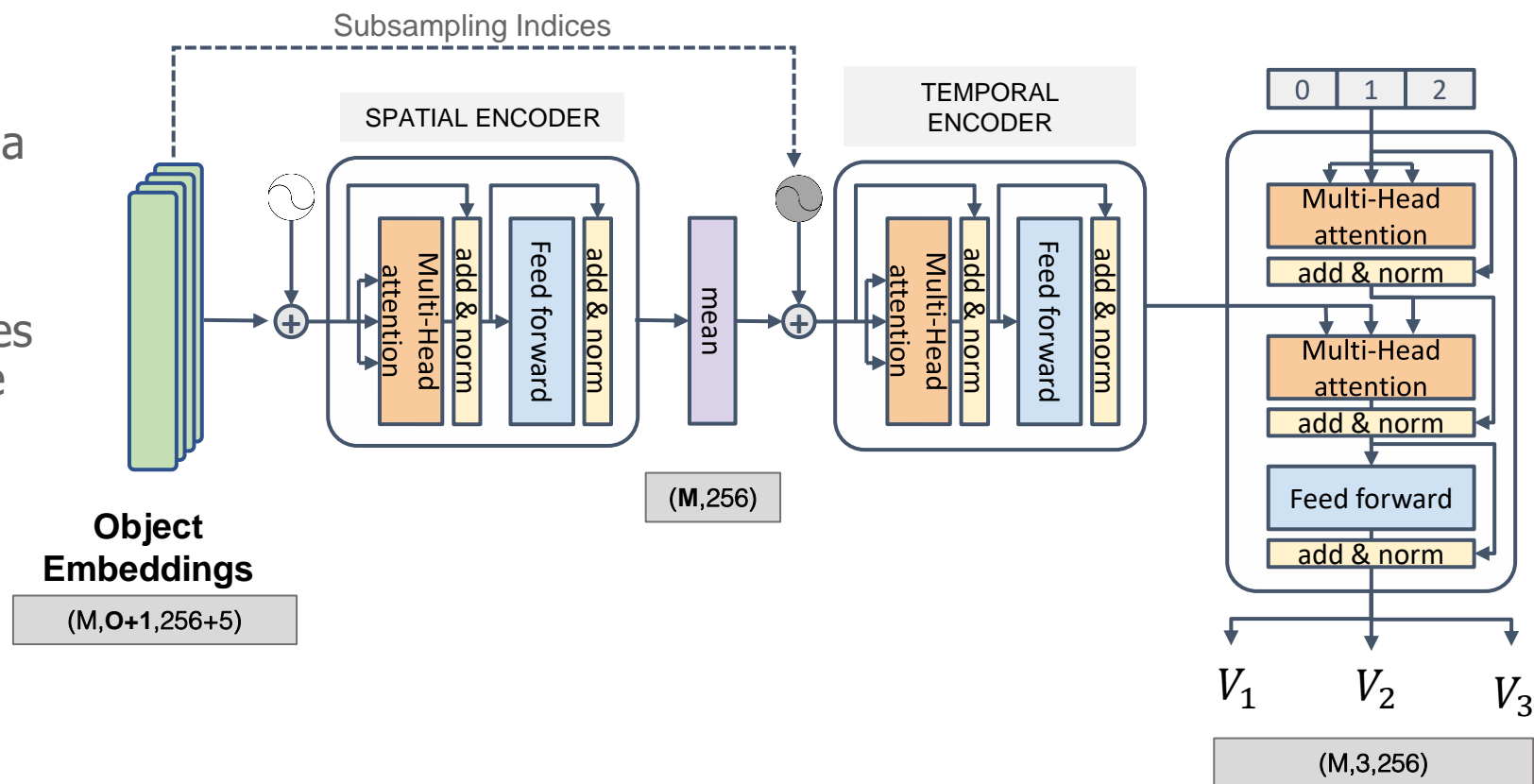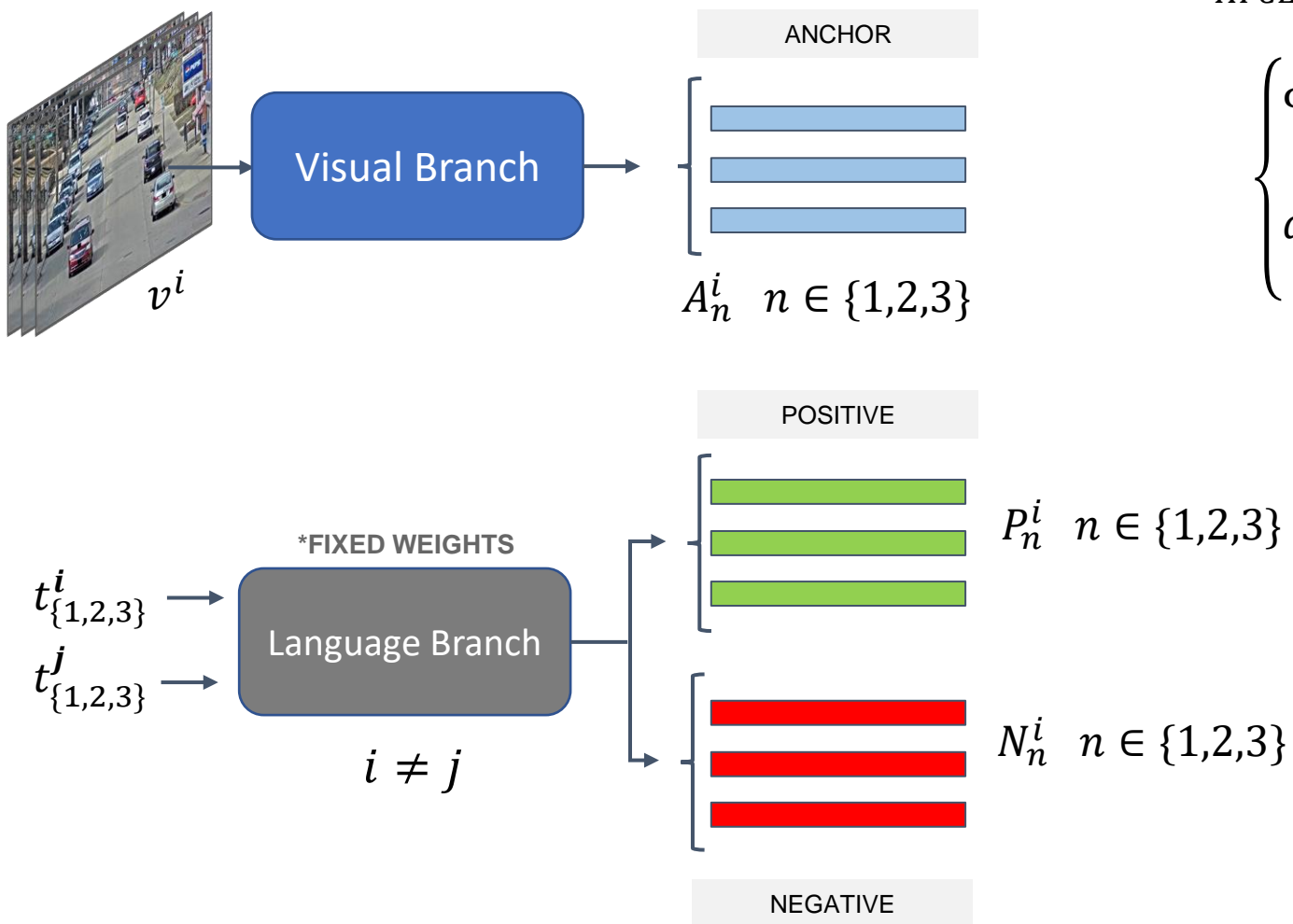[2]: Faster r-cnn: Towards real-time object detection with region proposal networks

# Spatio-Temporal Transformers

**RATIONALE:**

- The Spatial Encoder provides a comprehensive frame-level representation

- The Temporal Encoder encodes the sequential meaning of the depicted scene.

# Optimization

ANCHOR

Visual Branch

$v^i$

$A_n^i \quad n \in \{1,2,3\}$

POSITIVE

*FIXED WEIGHTS

$t_{\{1,2,3\}}^i$

$t_{\{1,2,3\}}^j$

Language Branch

$i \neq j$

$P_n^i \quad n \in \{1,2,3\}$

$N_n^i \quad n \in \{1,2,3\}$

NEGATIVE

$$\mathcal{L}_{AYCE}(A, P, N) = TL(A, P, N) + \frac{1}{Bs}\sum_{i=1}^{Bs}\beta \cdot \Phi(A^i, P^i)$$

$$\begin{cases} \Phi(A, P) = \min\left(d(A_m, P_n)\right) \quad m, n \in \{1,2,3\} \\ d_{TL}(A, [P|N]) = \frac{1}{9}\sum_{m=1}^{3}\sum_{n=1}^{3}\|A_m - [P|N]_n\|^2 \end{cases}$$

## RESULTS

| MODEL | MRR | RE@5 | RE@10 |
|---|---|---|---|
| **BASELINE** | 0.0269 | 0.0264 | 0.0491 |
| **OURS** (BEST − 11° ) | **0.1078** | 0.1321 | 0.2491 |
| Alibaba-UTS-ZJU (*1°*) | 0.1869 | - | - |
| Sun Asterisk (4°) | 0.1571 | - | - |
| Modulab (10° ) | 0. 1195 | - | - |

# Qualitative Results

"A red pickup drives straight down a highway."

"A red pickup truck runs down the street."

"A red pickup following straight other three car"



"A blue pickup runs down the street."

"A blue pickup truck going straight down the street passing an intersection."

"A blue truck runs down the street."



"A black sedan crossing an intersection."

"A midsize black sedan goes straight through the intersection."

"A black Sedan runs down the street."