

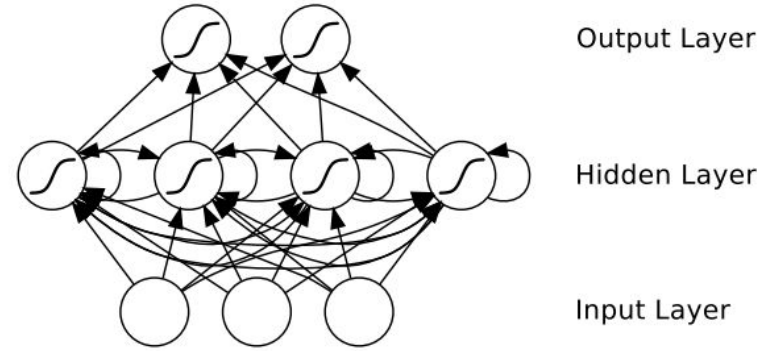
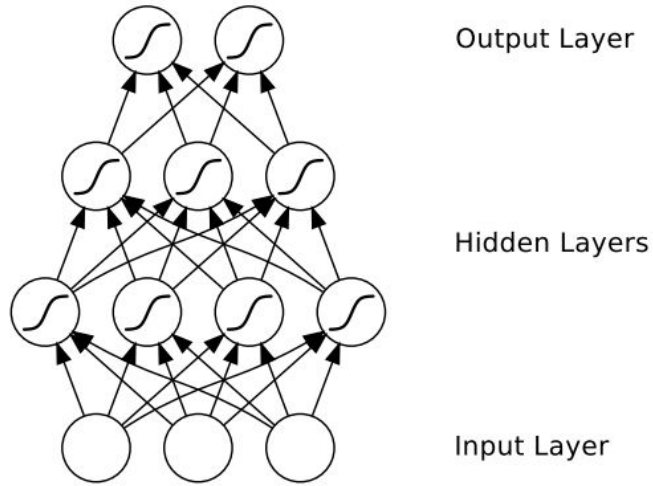
[Link](#)

RNN & LSTM

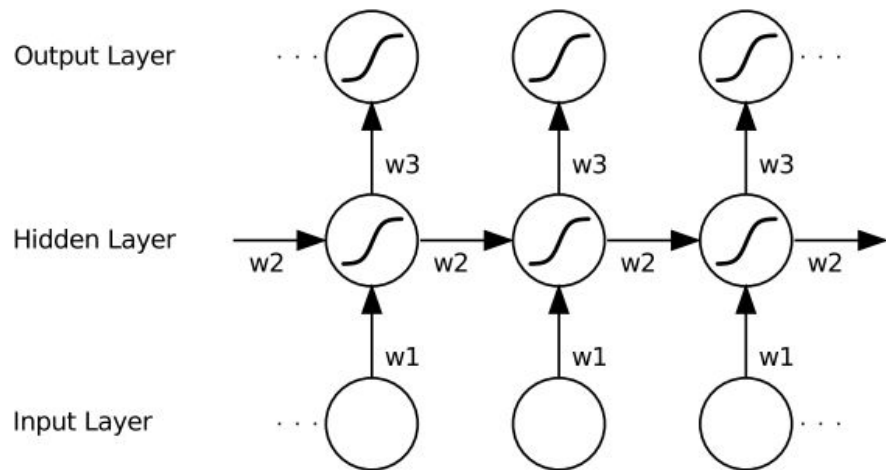
(Memory Networks)

Robotics Summer School
Sanket Kalwar
(31/05/2025)
Day 11

Neural Networks & RNN's

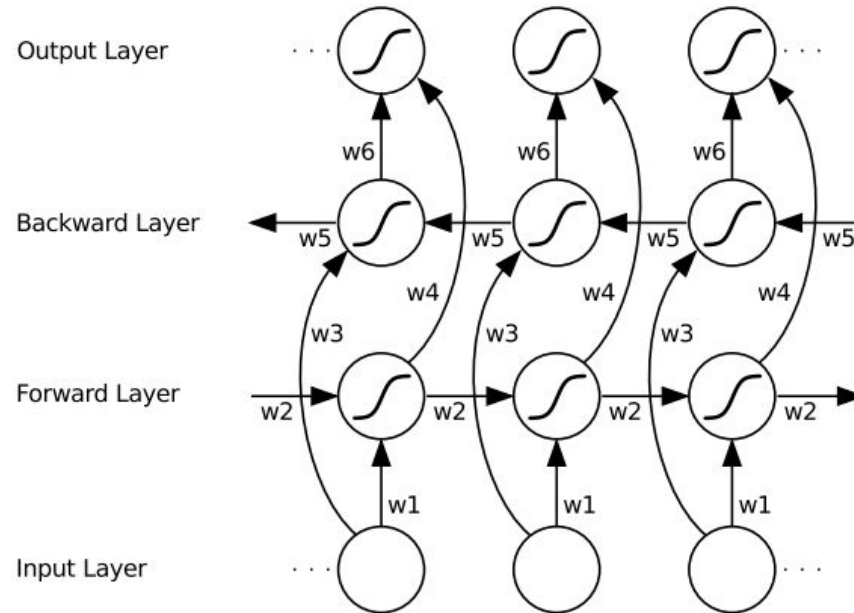


RNN:

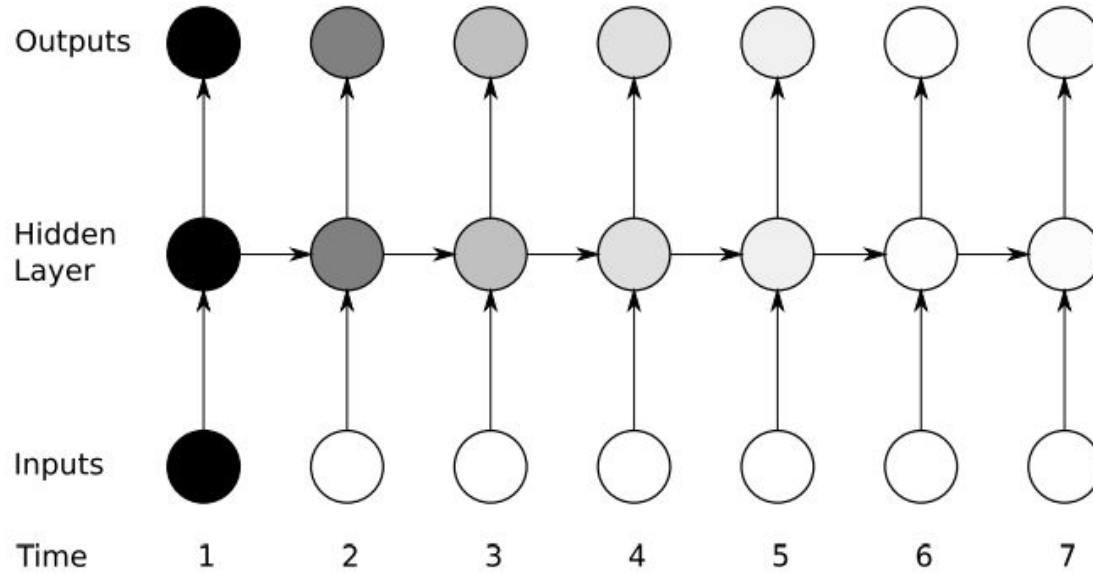


$$\begin{aligned}i_t &= \sigma(W_{ii}x_t + b_{ii} + W_{hi}h_{t-1} + b_{hi}) \\f_t &= \sigma(W_{if}x_t + b_{if} + W_{hf}h_{t-1} + b_{hf}) \\g_t &= \tanh(W_{ig}x_t + b_{ig} + W_{hg}h_{t-1} + b_{hg}) \\o_t &= \sigma(W_{io}x_t + b_{io} + W_{ho}h_{t-1} + b_{ho}) \\c_t &= f_t \odot c_{t-1} + i_t \odot g_t \\h_t &= o_t \odot \tanh(c_t)\end{aligned}$$

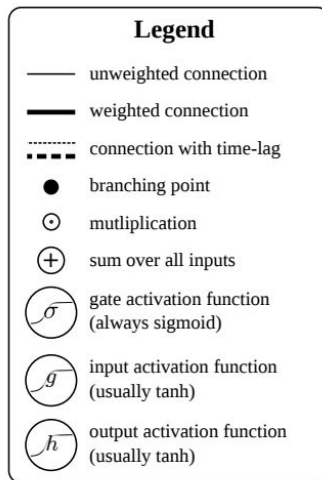
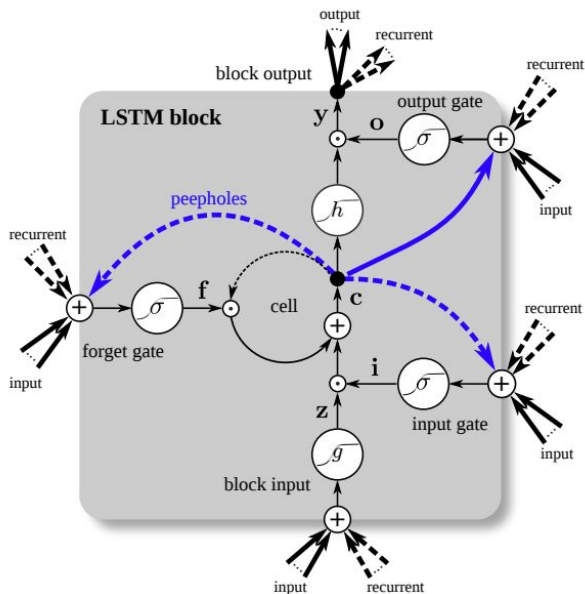
Bidirectional RNN:



Vanishing gradient Issue with RNN:



LSTM Architecture:



$$\bar{z}^t = W_z x^t + R_z y^{t-1} + b_z$$

$$z^t = g(\bar{z}^t)$$

$$\bar{i}^t = W_i x^t + R_i y^{t-1} + p_i \odot c^{t-1} + b_i$$

$$i^t = \sigma(\bar{i}^t)$$

$$\bar{f}^t = W_f x^t + R_f y^{t-1} + p_f \odot c^{t-1} + b_f$$

$$f^t = \sigma(\bar{f}^t)$$

$$c^t = z^t \odot i^t + c^{t-1} \odot f^t$$

$$\bar{o}^t = W_o x^t + R_o y^{t-1} + p_o \odot c^t + b_o$$

$$o^t = \sigma(\bar{o}^t)$$

$$y^t = h(c^t) \odot o^t$$

block input

input gate

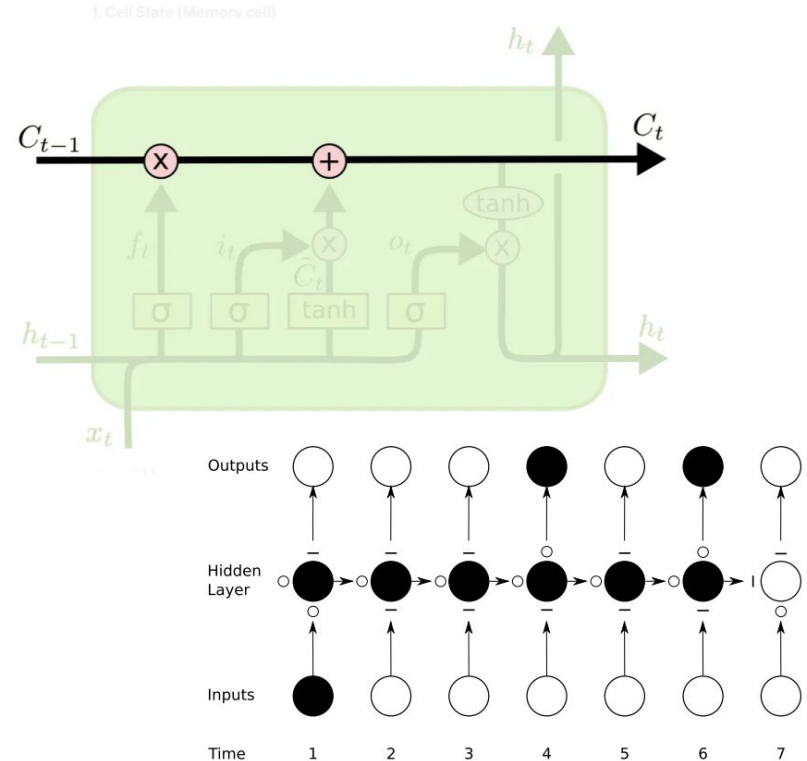
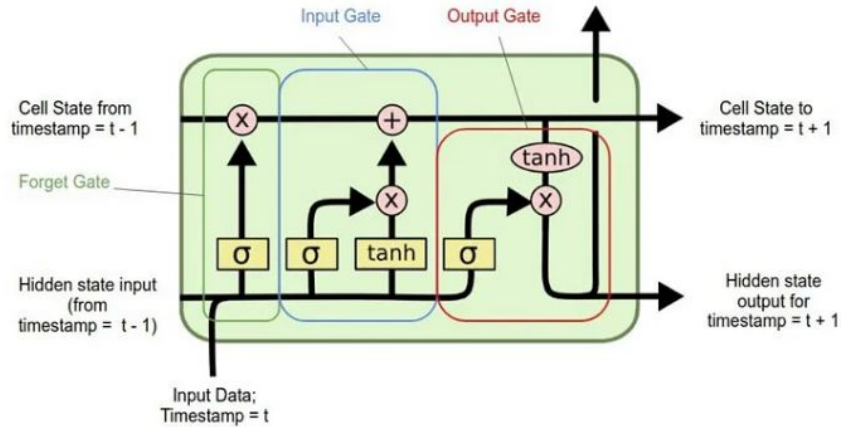
forget gate

cell

output gate

block output

Skip Connections (Addressing Vanishing Gradient Issue):



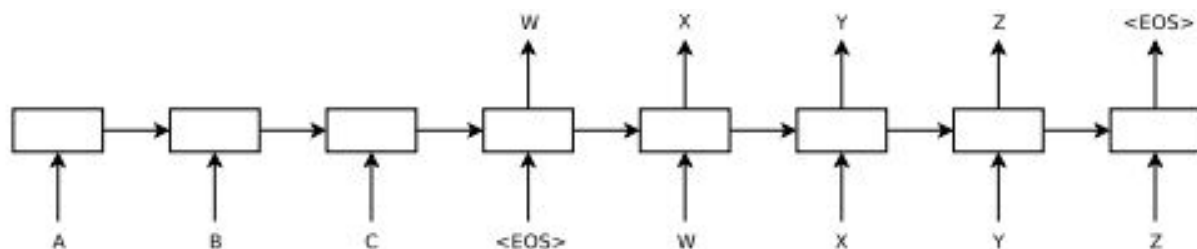
Seq2Seq Model:

Sequence to Sequence Learning with Neural Networks

Ilya Sutskever
Google
ilyasu@google.com

Oriol Vinyals
Google
vinyals@google.com

Quoc V. Le
Google
qvl@google.com



[Paper](#)

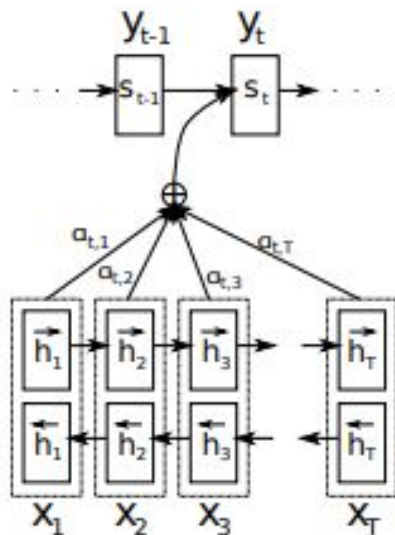
[Code](#)

NEURAL MACHINE TRANSLATION BY JOINTLY LEARNING TO ALIGN AND TRANSLATE

NEURAL MACHINE TRANSLATION BY JOINTLY LEARNING TO ALIGN AND TRANSLATE

Dzmitry Bahdanau
Jacobs University Bremen, Germany

KyungHyun Cho **Yoshua Bengio***
Université de Montréal



ABSTRACT

Neural machine translation is a recently proposed approach to machine translation. Unlike the traditional statistical machine translation, the neural machine translation aims at building a single neural network that can be jointly tuned to maximize the translation performance. The models proposed recently for neural machine translation often belong to a family of encoder-decoders and encode a source sentence into a fixed-length vector from which a decoder generates a translation. In this paper, we conjecture that the use of a fixed-length vector is a bottleneck in improving the performance of this basic encoder-decoder architecture, and propose to extend this by allowing a model to automatically (soft-)search for parts of a source sentence that are relevant to predicting a target word, without having to form these parts as a hard segment explicitly. With this new approach, we achieve a translation performance comparable to the existing state-of-the-art phrase-based system on the task of English-to-French translation. Furthermore, qualitative analysis reveals that the (soft-)alignments found by the model agree well with our intuition.

Show, Attend and Tell: Neural Image Caption Generation with Visual Attention

Show, Attend and Tell: Neural Image Caption Generation with Visual Attention

Kelvin Xu
Jimmy Lei Ba
Ryan Kiros
Kyunghyun Cho
Aaron Courville
Ruslan Salakhutdinov
Richard S. Zemel
Yoshua Bengio



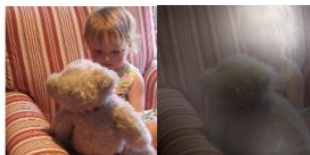
A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.



A stop sign is on a road with a mountain in the background.



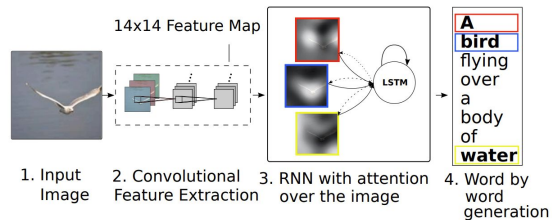
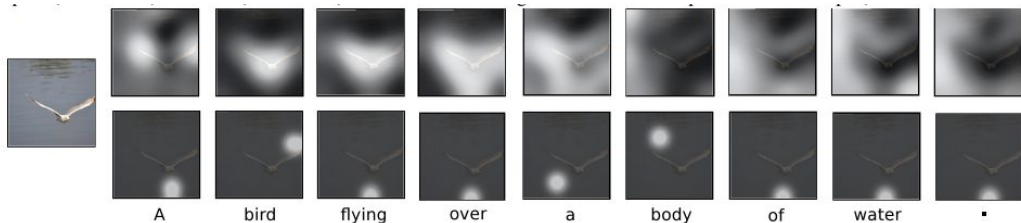
A little girl sitting on a bed with a teddy bear.



A group of people sitting on a boat in the water.



A giraffe standing in a forest with trees in the background.



[Paper](#)

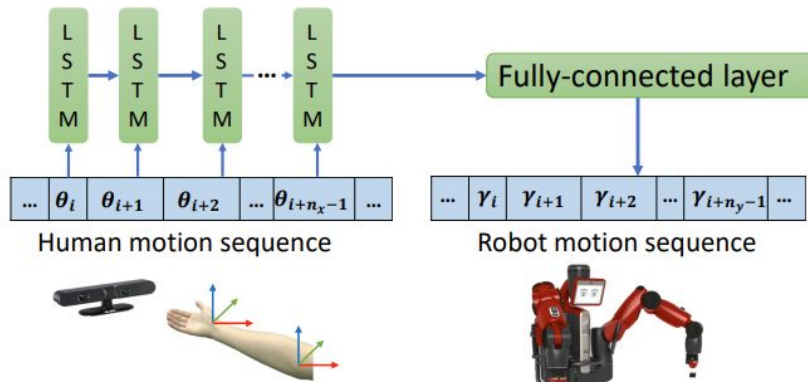
[Code](#)

Application To Robotics:

Collaborative Human-Robot Motion Generation using LSTM-RNN

Xuan Zhao, Sakmongkon Chumkamon, Shuanda Duan, Juan Rojas, and Jia Pan[†]

Abstract—We propose a deep learning based method for fast and responsive human-robot handovers that generate robot motion according to human motion observations. Our method learns an offline human-robot interaction model through a Recurrent Neural Network with Long Short-Term Memory units (LSTM-RNN). The robot uses the learned network to respond appropriately to novel online human motions. Our method is tested both on pre-recorded data and real-world human-robot handover experiments. Our method achieves robot motion accuracies that outperform the baseline. In addition, our method demonstrates a strong ability to adapt to changes in velocity of human motions.



Prior Work From RRC Using LSTM:

ATPPNet: Attention based Temporal Point cloud Prediction Network

Kaustab Pal^{*1}, Aditya Sharma^{*1}, Avinash Sharma², K. Madhava Krishna¹

Abstract—Point cloud prediction is an important yet challenging task in the field of autonomous driving. The goal is to predict future point cloud sequences that maintain object structures while accurately representing their temporal motion. These predicted point clouds help in other subsequent tasks like object trajectory estimation for collision avoidance or estimating locations with the least odometry drift. In this work, we present ATPPNet, a novel architecture that predicts future point cloud sequences given a sequence of previous time step point clouds obtained with LiDAR sensor. ATPPNet leverages Conv-LSTM along with channel-wise and spatial attention dually complemented by a 3D-CNN branch for extracting an enhanced spatio-temporal context to recover high quality fidel predictions of future point clouds. We conduct extensive experiments on publicly available datasets and report impressive performance outperforming the existing methods. We also conduct a thorough ablative study of the proposed architecture and provide an application study that highlights the potential of our model for tasks like odometry estimation.

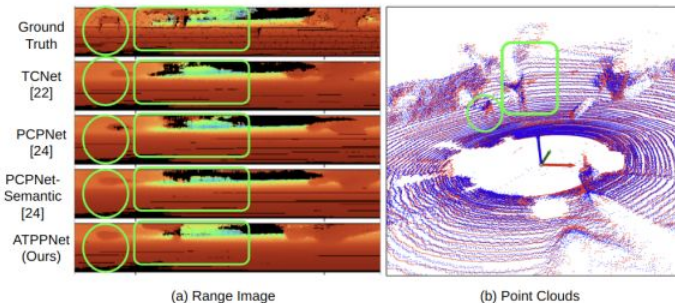
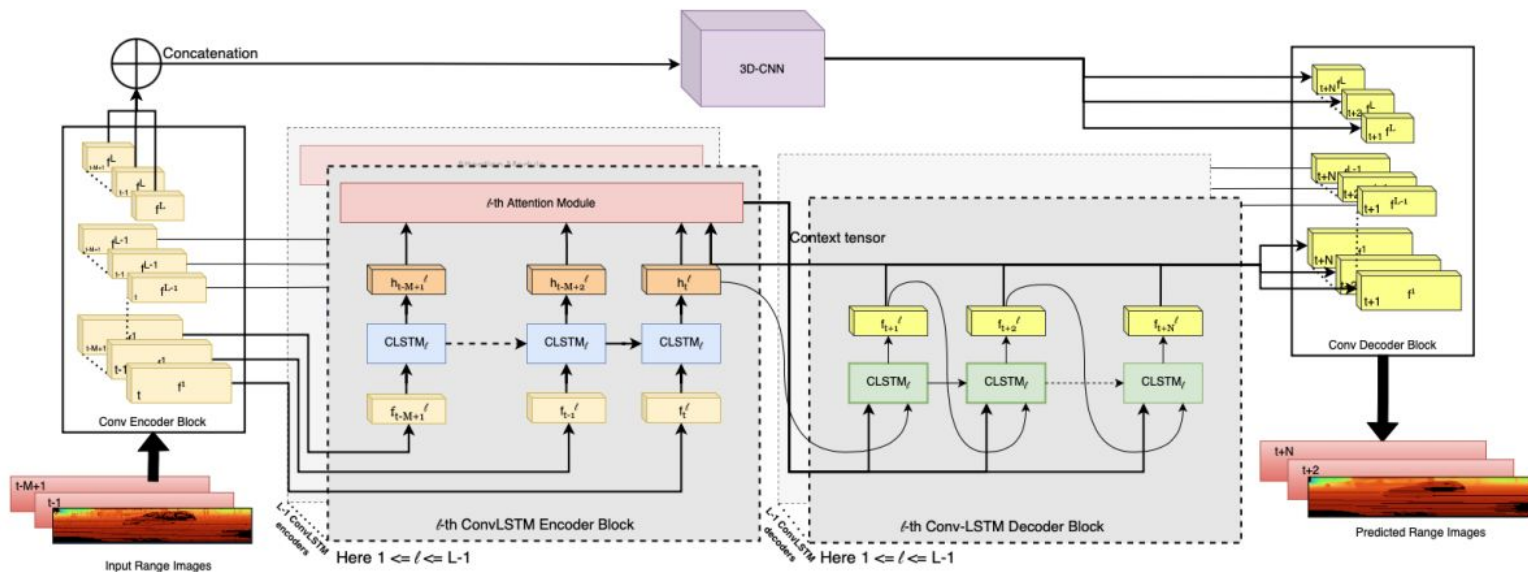


Fig. 1: (a) Predicted range images by our ATPPNet and existing methods in comparison to ground truth and, (b) the 3D rendering of the predicted point cloud by ATPPNet (blue) and ground-truth (red). Green circle/rectangle highlights regions where ATPPNet's predictions are superior.

(e.g., CNNs) and sequence prediction (e.g., LSTMs) cannot be directly employed as they cannot process spatially unordered data. Another key challenge is that the LiDAR point clouds are extremely sparse making it difficult to capture

ATPPNet: Attention based Temporal Point cloud Prediction Network



Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

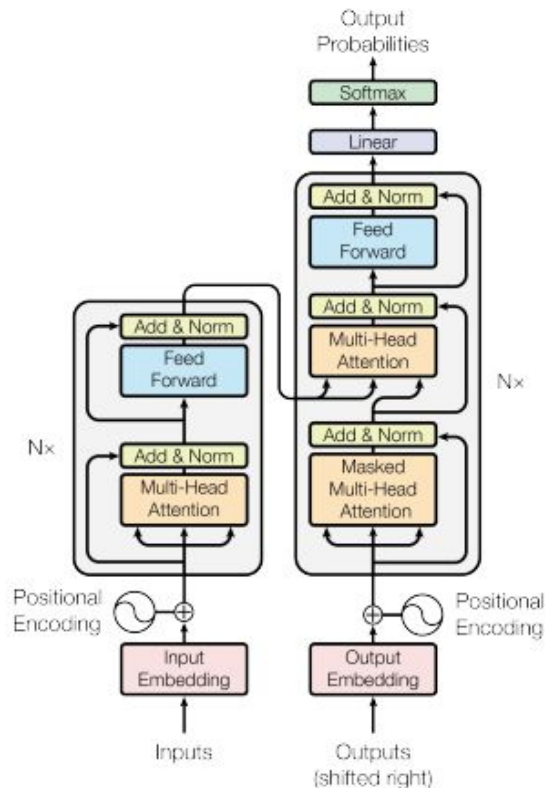
Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Lukasz Kaiser*
Google Brain
lukaszkaizer@google.com

Illia Polosukhin* ‡
illia.polosukhin@gmail.com

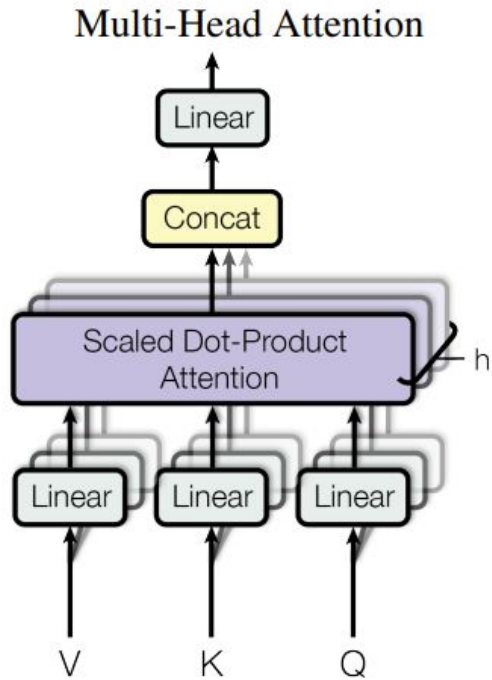
Transformer (Memory Networks)

Attention is ALL You Need

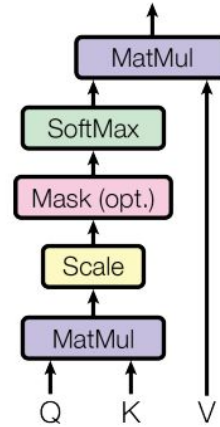


[Paper](#) [Code](#)

MHA Layer:



Scaled Dot-Product Attention



Need Some Attention!

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

$$\text{where head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

Where the projections are parameter matrices $W_i^Q \in \mathbb{R}^{d_{\text{model}} \times d_k}$, $W_i^K \in \mathbb{R}^{d_{\text{model}} \times d_k}$, $W_i^V \in \mathbb{R}^{d_{\text{model}} \times d_v}$ and $W^O \in \mathbb{R}^{hd_v \times d_{\text{model}}}$

$$\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2$$

$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$

$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$

Let's Talk about Compute:

Layer Type	Complexity per Layer	Sequential Operations	Maximum Path Length
Self-Attention	$O(n^2 \cdot d)$	$O(1)$	$O(1)$
Recurrent	$O(n \cdot d^2)$	$O(n)$	$O(n)$
Convolutional	$O(k \cdot n \cdot d^2)$	$O(1)$	$O(\log_k(n))$
Self-Attention (restricted)	$O(r \cdot n \cdot d)$	$O(1)$	$O(n/r)$

1st Attempt To Make transformer work on Images:

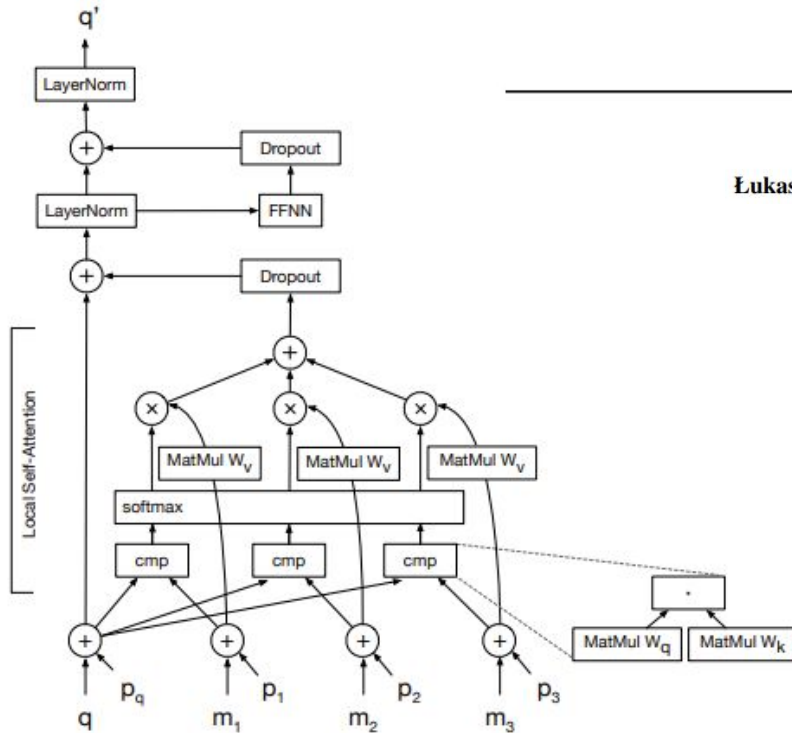


Image Transformer

Niki Parmar ^{*1} **Ashish Vaswani** ^{*1} **Jakob Uszkoreit** ¹
Łukasz Kaiser ¹ **Noam Shazeer** ¹ **Alexander Ku** ^{2,3} **Dustin Tran** ⁴

2nd Attempt To Make transformer work on Images:

Non-local Neural Networks

Xiaolong Wang^{1,2*}

Ross Girshick²

Abhinav Gupta¹

Kaiming He²

¹Carnegie Mellon University

²Facebook AI Research

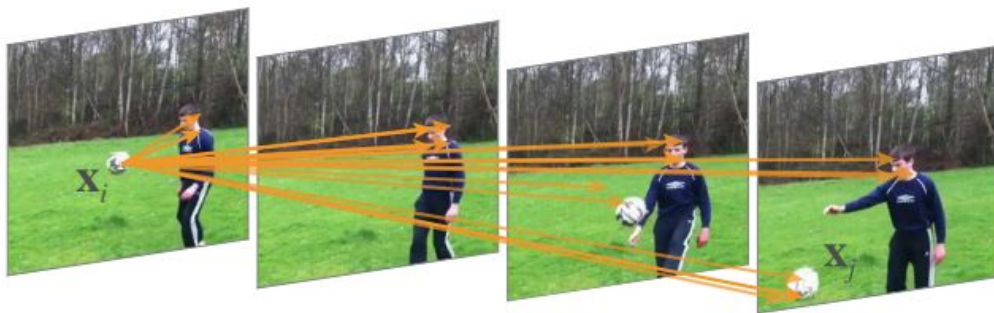
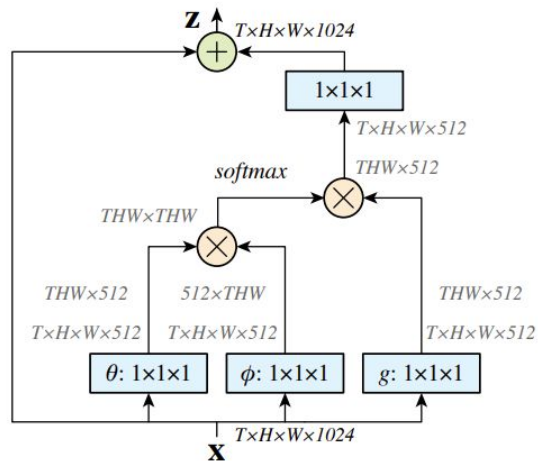


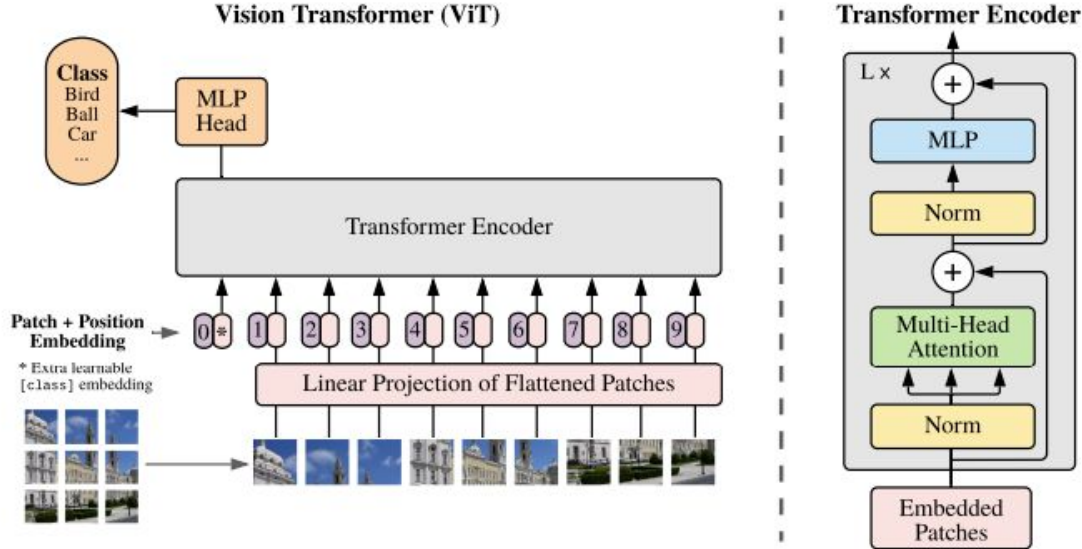
Figure 1. A spacetime *non-local* operation in our network trained for video classification in Kinetics. A position \mathbf{x}_i 's response is computed by the weighted average of the features of *all* positions \mathbf{x}_j (only the highest weighted ones are shown here). In this example computed by our model, note how it relates the ball in the first frame to the ball in the last two frames. More examples are in Figure 3.



[Paper](#)

[Code](#)

AN IMAGE IS WORTH 16X16 WORDS: TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE



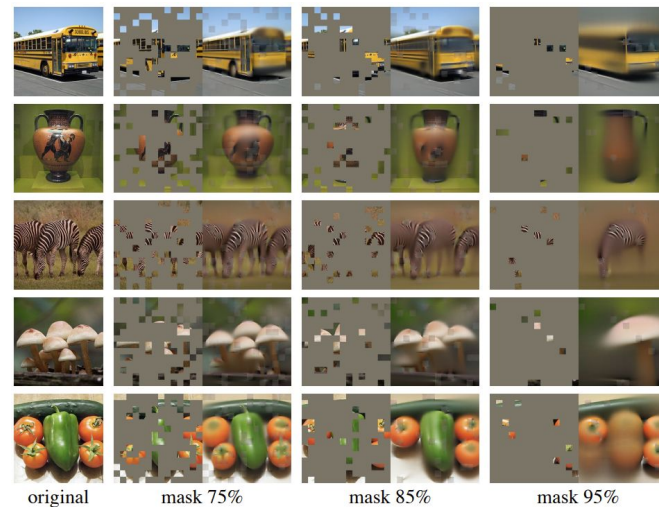
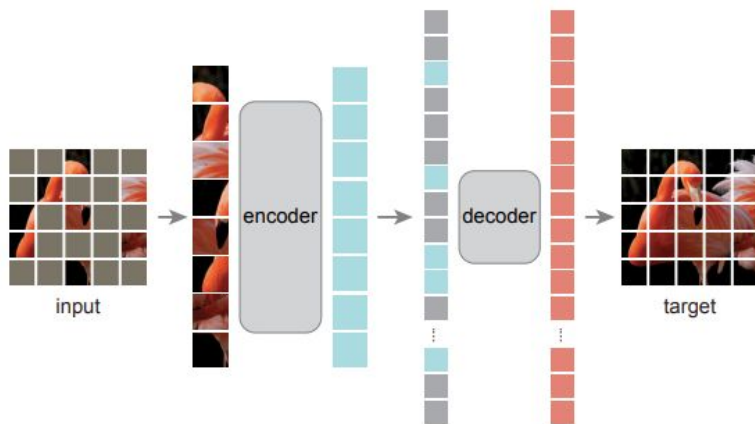
Masked Autoencoder (MAE):

Masked Autoencoders Are Scalable Vision Learners

Kaiming He^{*,†} Xinlei Chen^{*} Saining Xie Yanghao Li Piotr Dollár Ross Girshick

^{*}equal technical contribution [†]project lead

Facebook AI Research (FAIR)



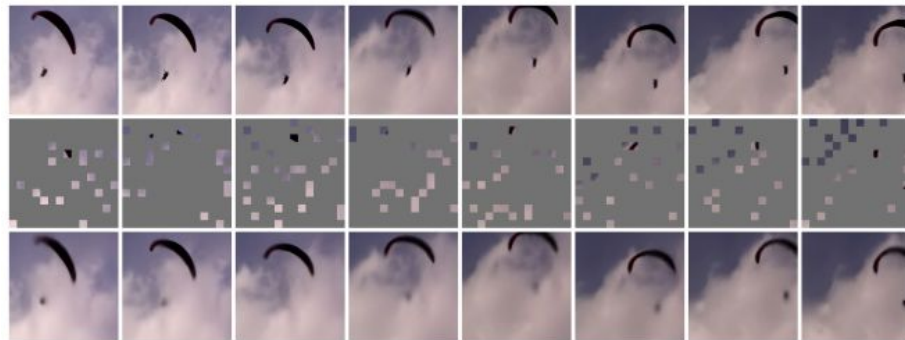
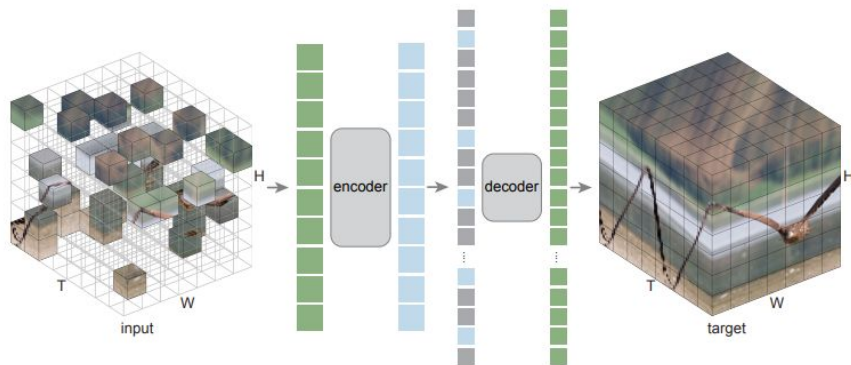
[Paper](#)

[Code](#)

Masked Autoencoders As Spatiotemporal Learners:

Masked Autoencoders As Spatiotemporal Learners

Christoph Feichtenhofer* Haoqi Fan* Yanghao Li Kaiming He
Meta AI, FAIR



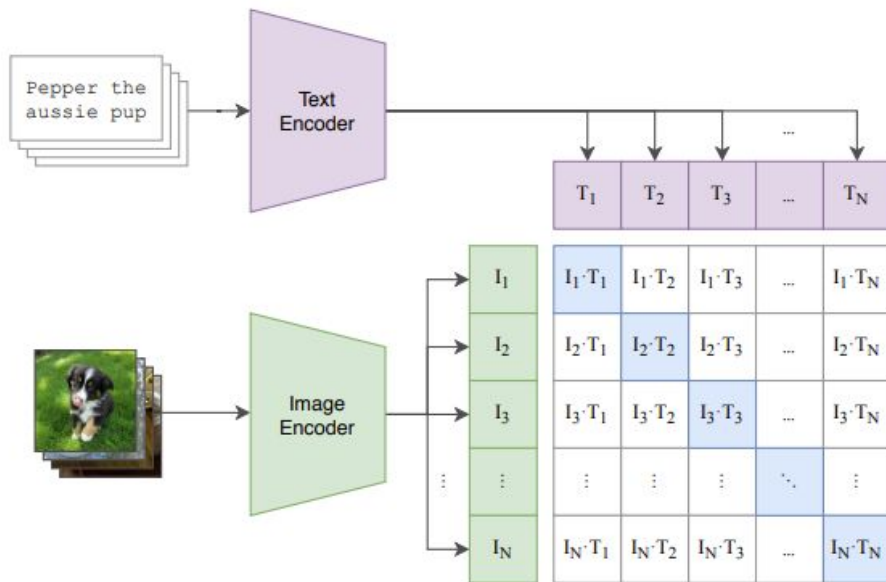
[Paper](#)

[Code](#)

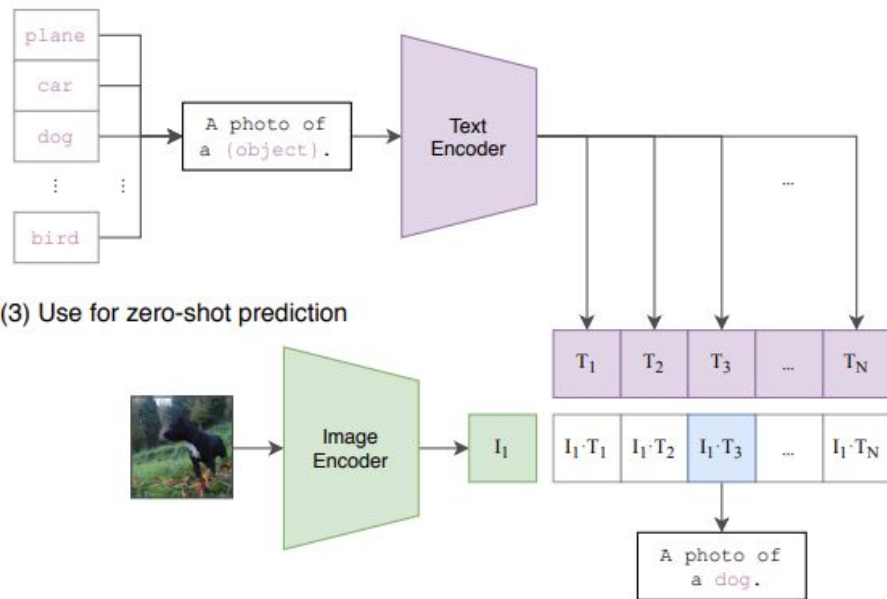
Visual-Language Model's

CLIP:

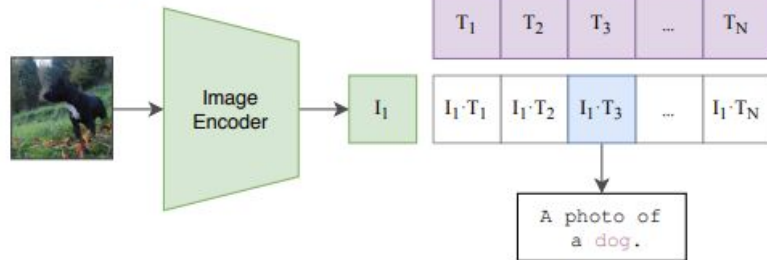
(1) Contrastive pre-training



(2) Create dataset classifier from label text



(3) Use for zero-shot prediction



[Paper](#)

[Code](#)

Training Pseudo code:

```
# image_encoder - ResNet or Vision Transformer
# text_encoder  - CBOW or Text Transformer
# I[n, h, w, c] - minibatch of aligned images
# T[n, l]       - minibatch of aligned texts
# W_i[d_i, d_e] - learned proj of image to embed
# W_t[d_t, d_e] - learned proj of text to embed
# t            - learned temperature parameter

# extract feature representations of each modality
I_f = image_encoder(I) #[n, d_i]
T_f = text_encoder(T)  #[n, d_t]

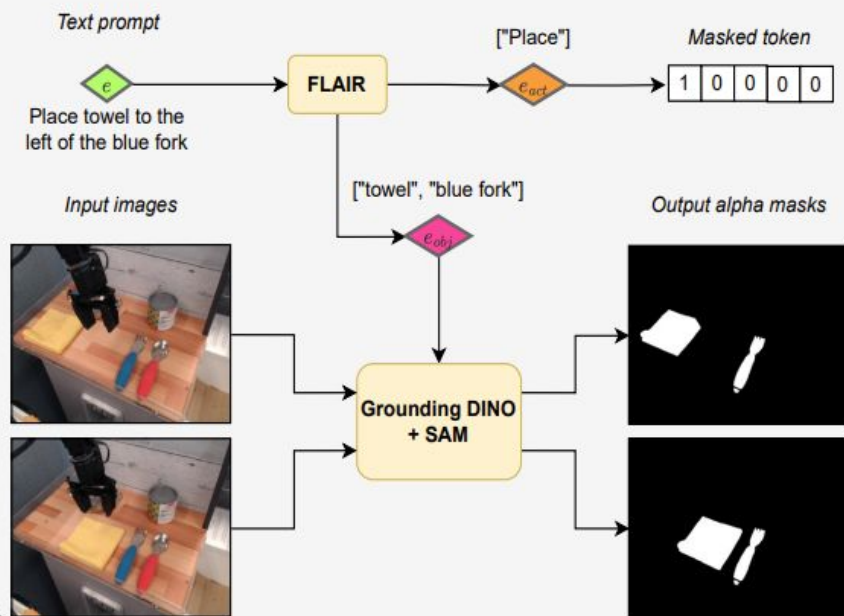
# joint multimodal embedding [n, d_e]
I_e = l2_normalize(np.dot(I_f, W_i), axis=1)
T_e = l2_normalize(np.dot(T_f, W_t), axis=1)

# scaled pairwise cosine similarities [n, n]
logits = np.dot(I_e, T_e.T) * np.exp(t)

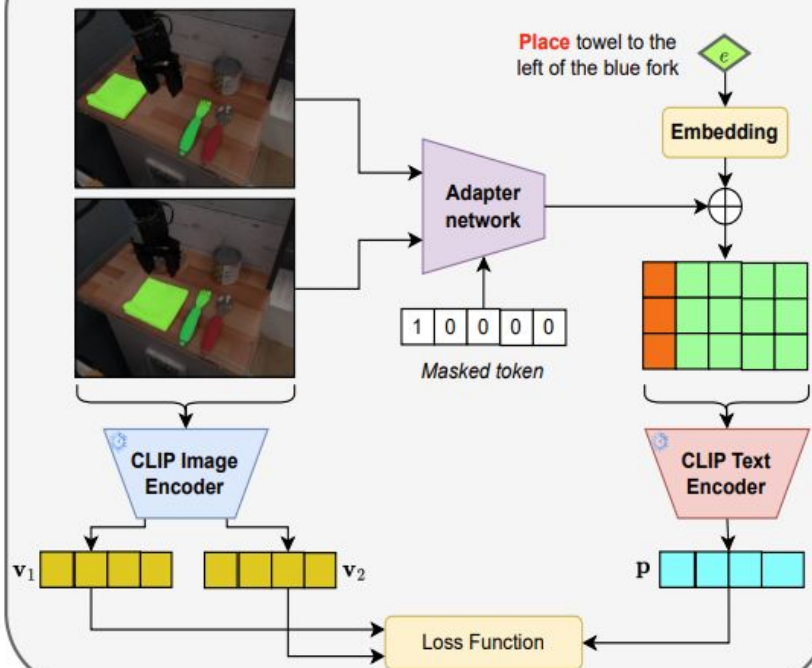
# symmetric loss function
labels = np.arange(n)
loss_i = cross_entropy_loss(logits, labels, axis=0)
loss_t = cross_entropy_loss(logits, labels, axis=1)
loss   = (loss_i + loss_t)/2
```

Robot-CLIP:

Dataset Preparation Pipeline



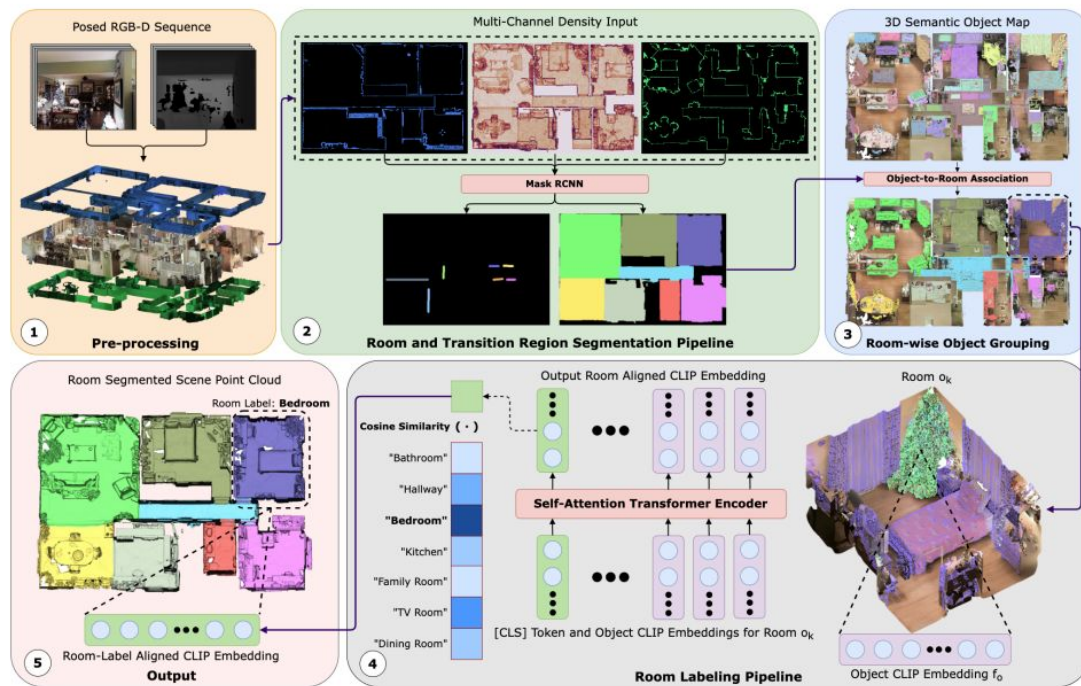
Fine-tuning Pipeline



Prior Work From RRC:

QueSTMaps: Queryable Semantic Topological Maps for 3D Scene Understanding

Yash Mehan^{1*}, Kumaraditya Gupta^{1*}, Rohit Jayanti^{1*}, Anirudh Govil¹, Sourav Garg², and Madhava Krishna¹

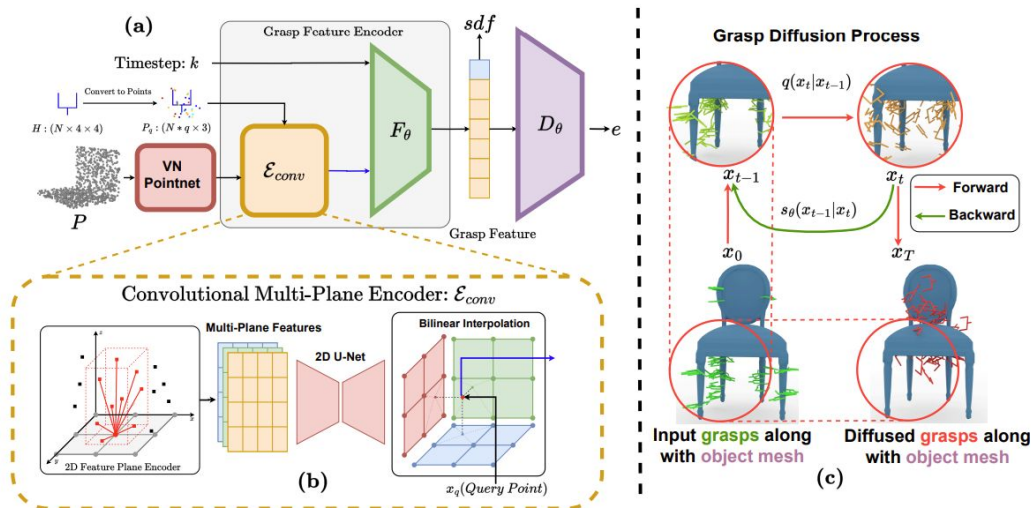
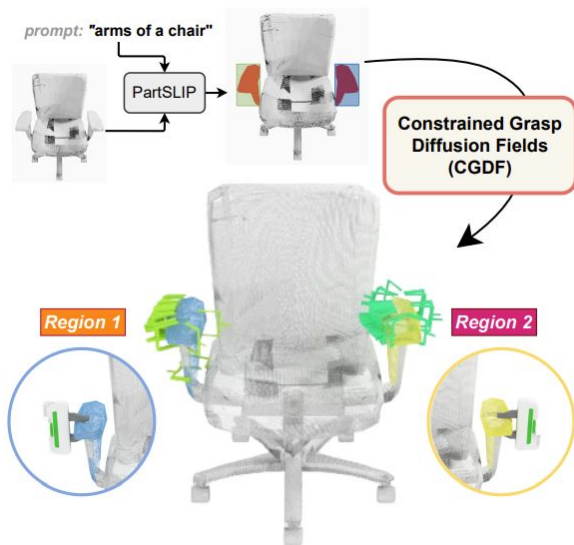


[Paper](#)

Prior Work From RRC:

Constrained 6-DoF Grasp Generation on Complex Shapes for Improved Dual-Arm Manipulation

Gaurav Singh^{1*}, Sanket Kalwar^{1*}, Md Faizal Karim¹, Bipasha Sen², Nagamanikandan Govindan¹,
Srinath Sridhar³ and K Madhava Krishna¹



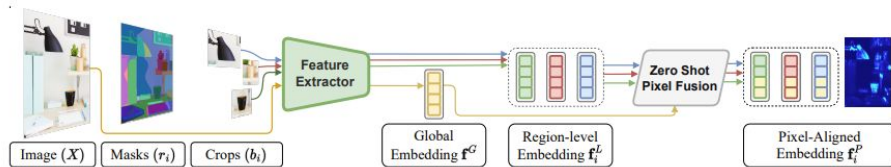
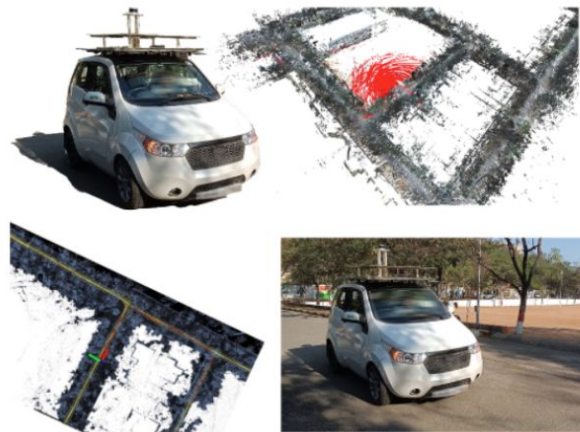
[Code](#) [Paper](#)

Prior Work From RRC:

ConceptFusion: Open-set Multimodal 3D Mapping

Krishna Murthy Jatavallabhula¹, Alihusein Kuwajerwala^{2,†}, Qiao Gu^{3,†}, Mohd Omama^{4,†}, Tao Chen¹, Alaa Maalouf¹, Shuang Li¹, Ganesh Iyer^{7,‡}, Soroush Saryazdi⁸, Nikhil Keetha⁵, Ayush Tewari¹, Joshua B. Tenenbaum¹, Celso Miguel de Melo⁶, K. Madhava Krishna⁴, Liam Paull², Florian Shkurti³, and Antonio Torralba¹

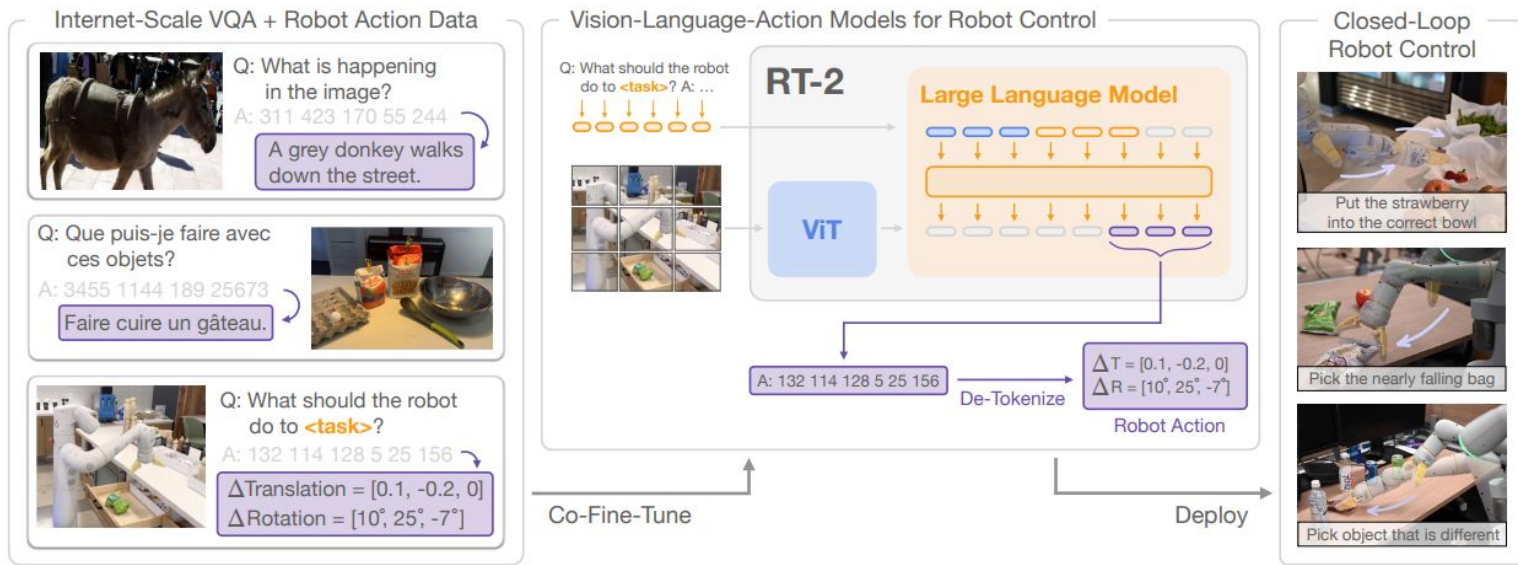
¹MIT, ²Université de Montréal, ³University of Toronto, ⁴IIIT Hyderabad, ⁵CMU, ⁶DEVCOM Army Research Lab, ⁷Amazon, ⁸Concordia University, [†]Co-second authors, [‡]Work done prior to Amazon



[Paper](#)

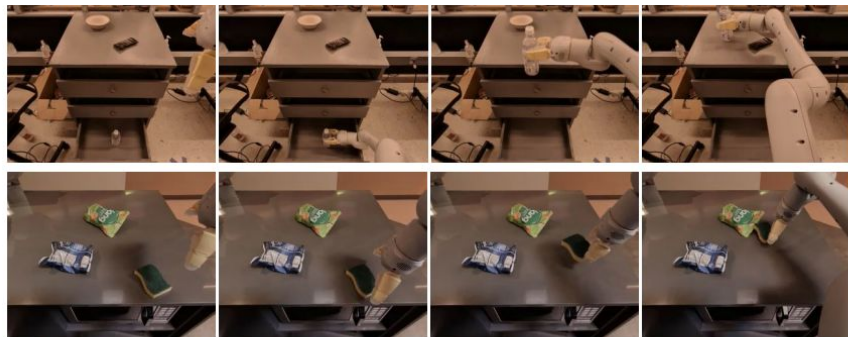
Application To Robotics:

RT-2: Vision-Language-Action Models Transfer Web Knowledge to Robotic Control



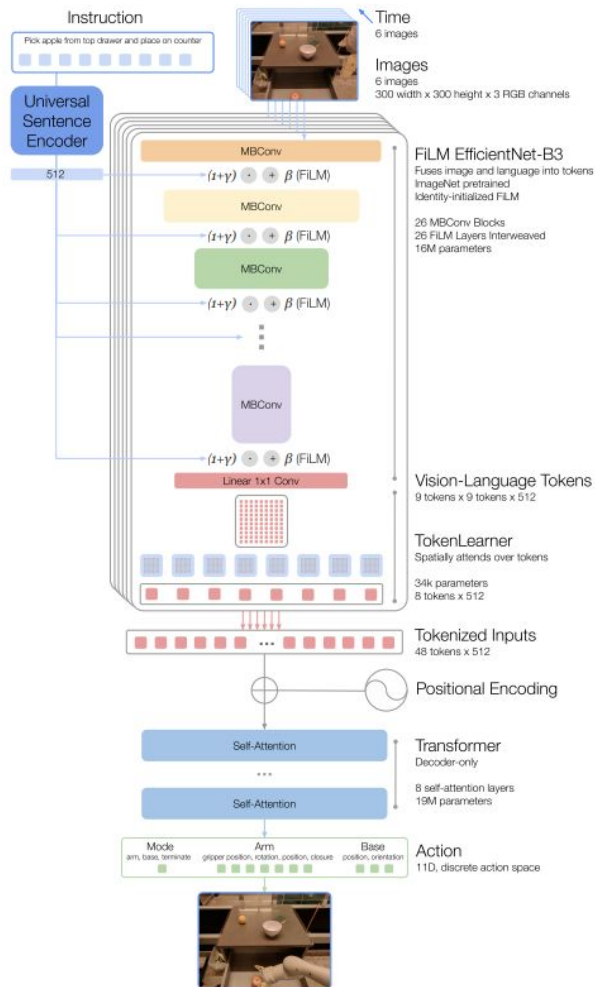
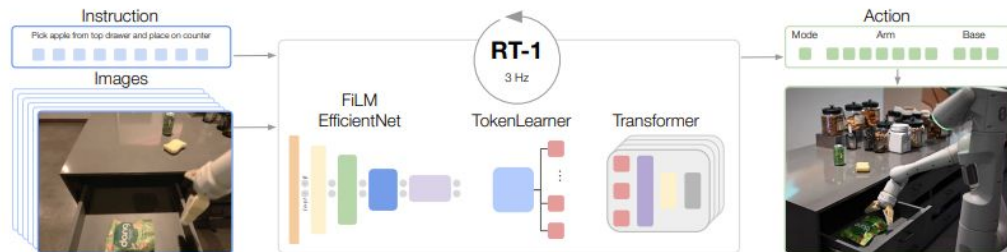
Application To Robotics:

RT-1: ROBOTICS TRANSFORMER FOR REAL-WORLD CONTROL AT SCALE



"pick water bottle from the bottom drawer and put it on the counter"

"move sponge to green jalapeno chips"



Application To Robotics:

Google DeepMind

Gemini Robotics: Bringing AI into the Physical World

Gemini Robotics Team, Google DeepMind¹

Dexterous, general & instructable Vision-Language-Action model



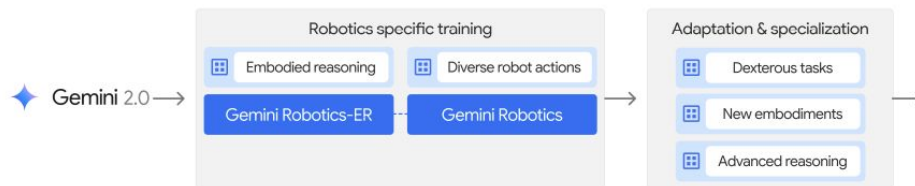
Complex dexterous tasks



New embodiments



Advanced reasoning & acting



Advanced embodied reasoning for robotics



DINO & DINOv2

DINO:

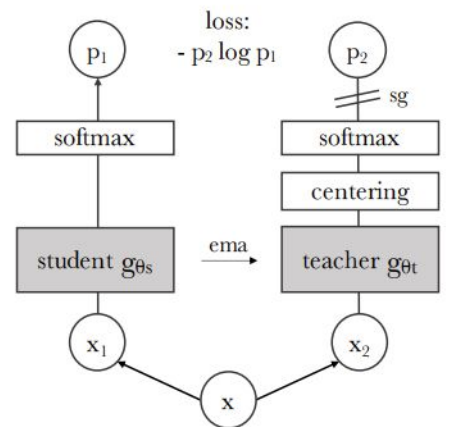
Supervised



DINO



	Random	Supervised	DINO
ViT-S/16	22.0	27.3	45.9
ViT-S/8	21.8	23.7	44.7

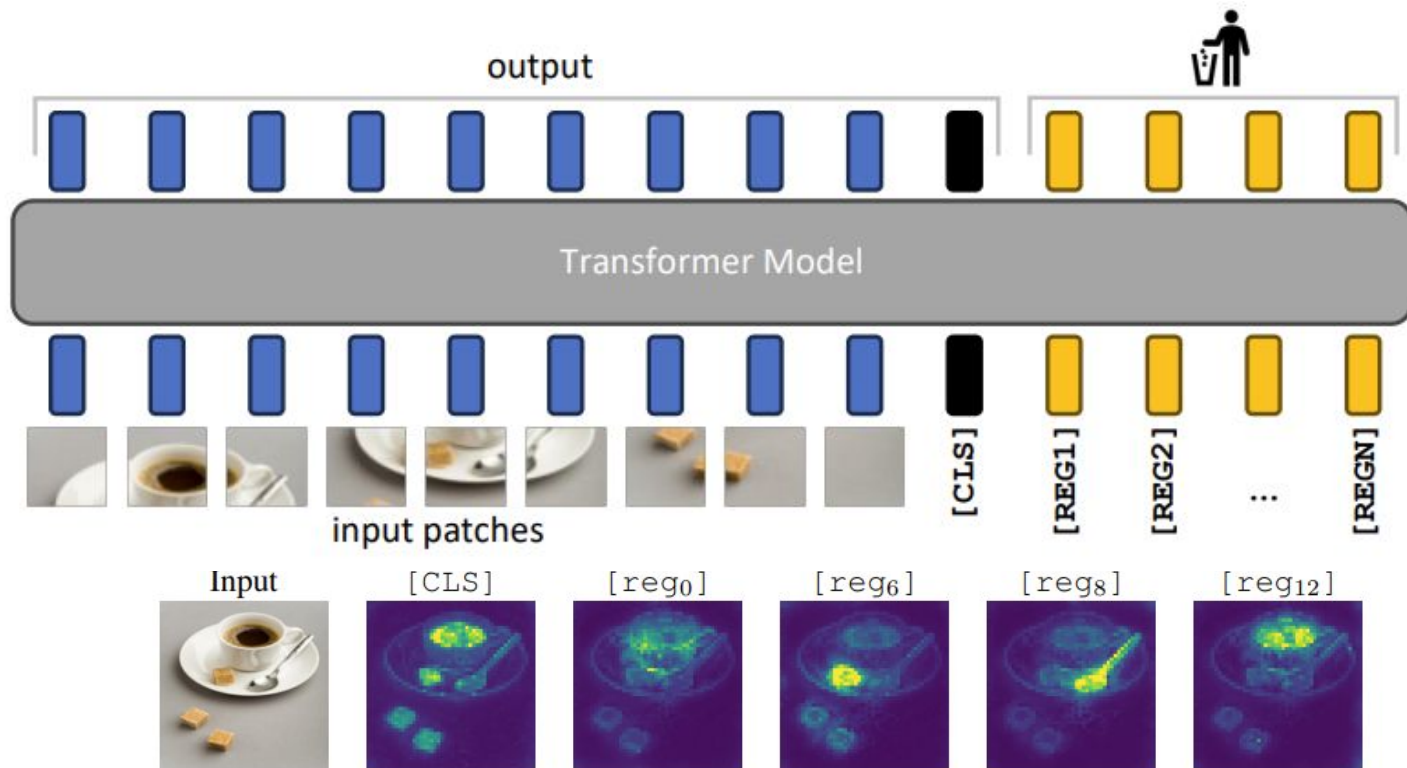


$$\min_{\theta_s} H(P_t(x), P_s(x))$$

[Code](#)

[Paper](#)

DINO2:



[Paper](#)

[Code](#)

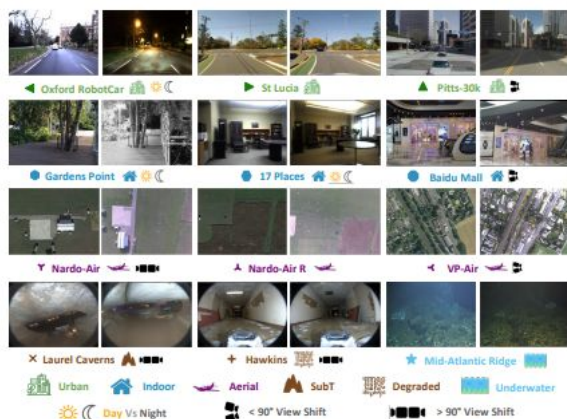
Prior Work From RRC:

AnyLoc: Towards Universal Visual Place Recognition

<https://anyloc.github.io/>

Nikhil Keetha^{*1}, Avneesh Mishra^{*2}, Jay Karhade^{*1}, Krishna Murthy Jatavallabhula³,
Sebastian Scherer¹, Madhava Krishna², and Sourav Garg⁴

¹CMU, ²IIT Hyderabad, ³MIT, ⁴University of Adelaide

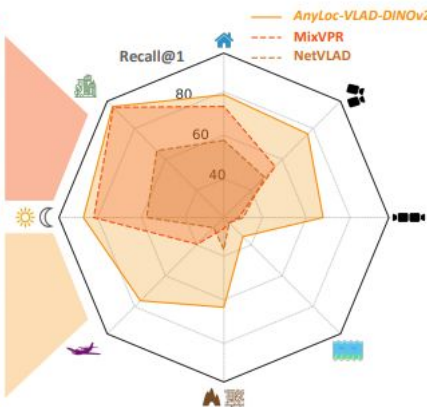
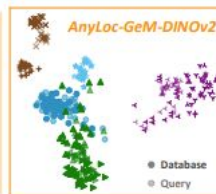
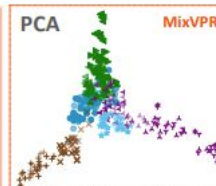


Large-Scale VPR Training

Supervised SOTA
VPR Baselines

Large-Scale Task-Agnostic
Pretraining \Rightarrow Freeze

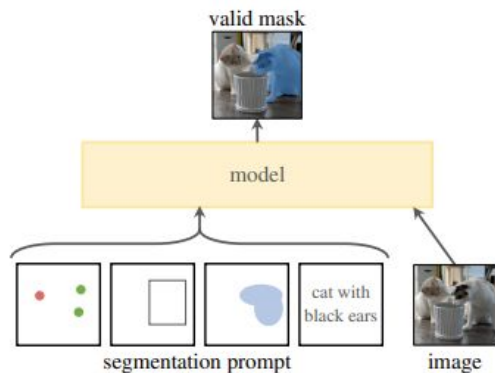
Foundation Model
Features



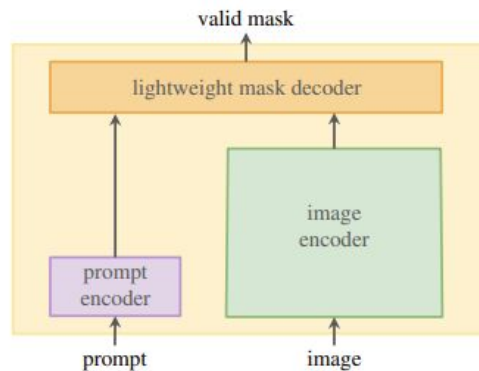
[Paper](#)

SAM & SAM2

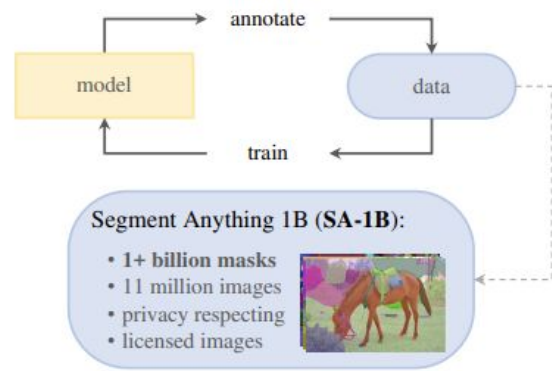
Segment Anything:



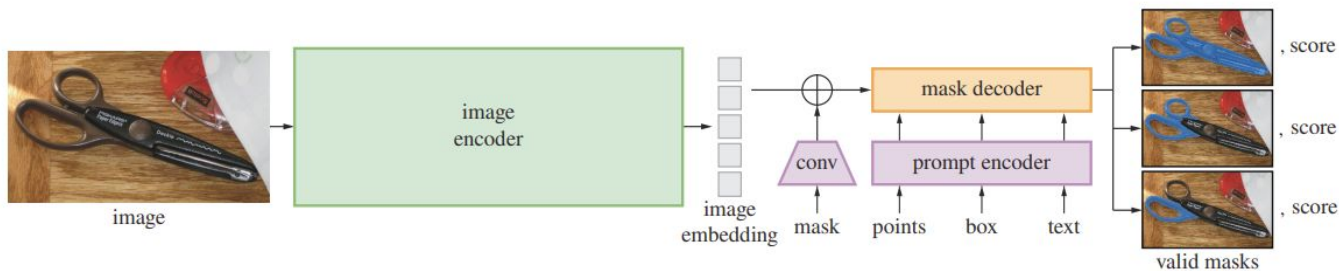
(a) **Task:** promptable segmentation



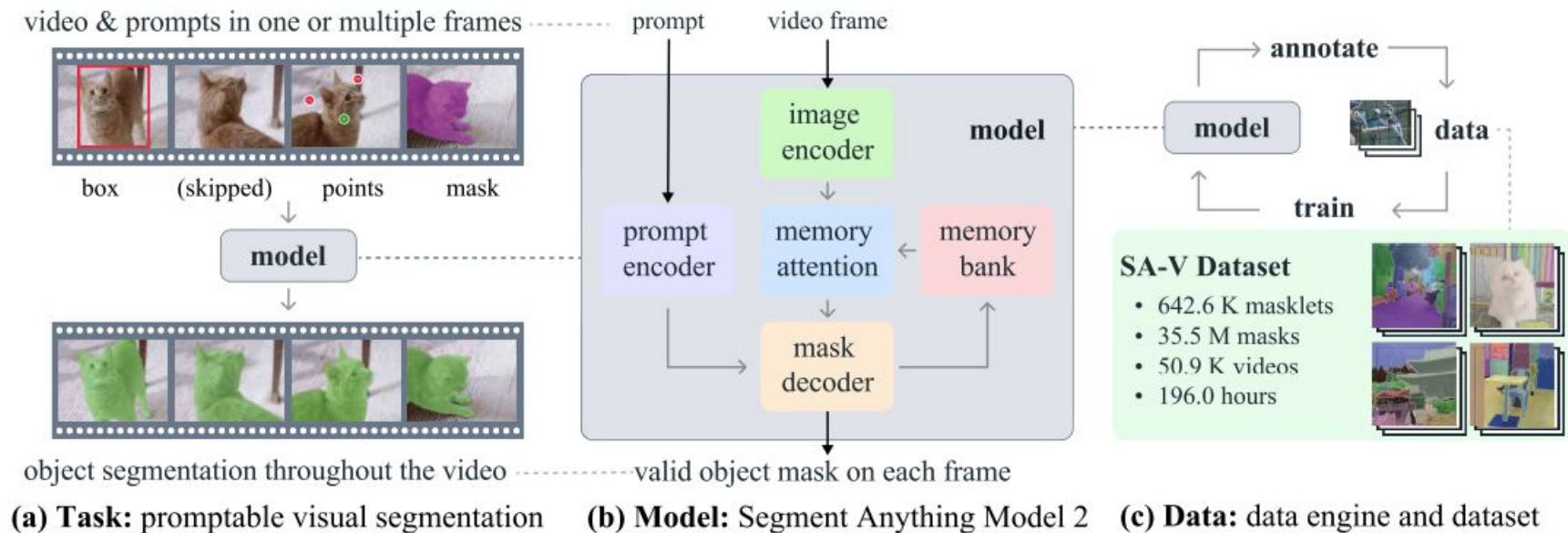
(b) **Model:** Segment Anything Model (SAM)



(c) **Data:** data engine (top) & dataset (bottom)



SAM2:



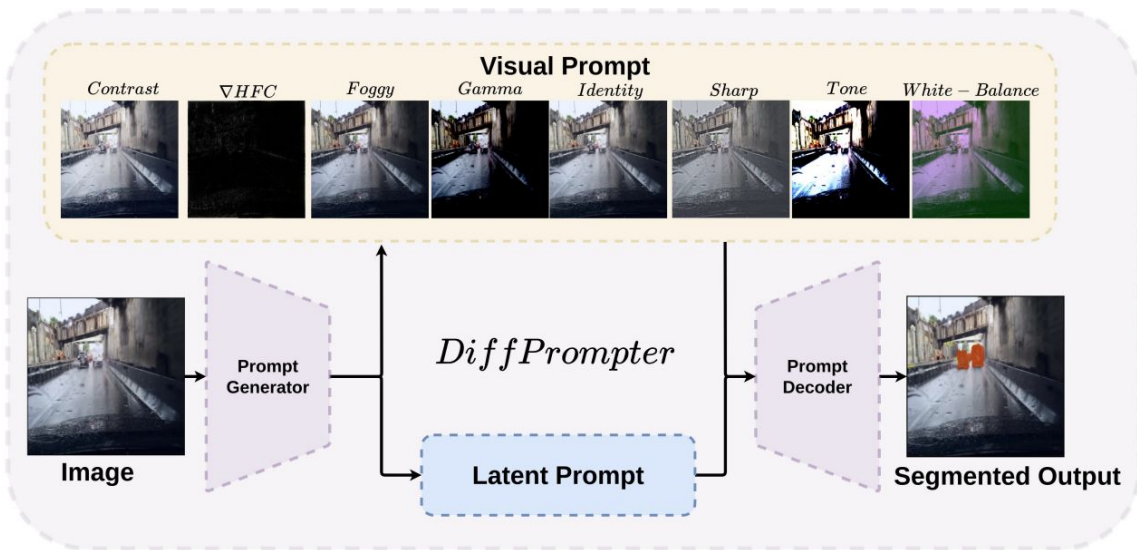
[Paper](#)

[Code](#)

Prior Work From RRC:

DiffPrompter: Differentiable Implicit Visual Prompts for Semantic-Segmentation in Adverse Conditions

Sanket Kalwar^{*1}, Mihir Ungarala^{*1}, Shruti Jain^{*1}, Aaron Monis¹, Krishna Reddy Konda³
Sourav Garg², K Madhava Krishna¹








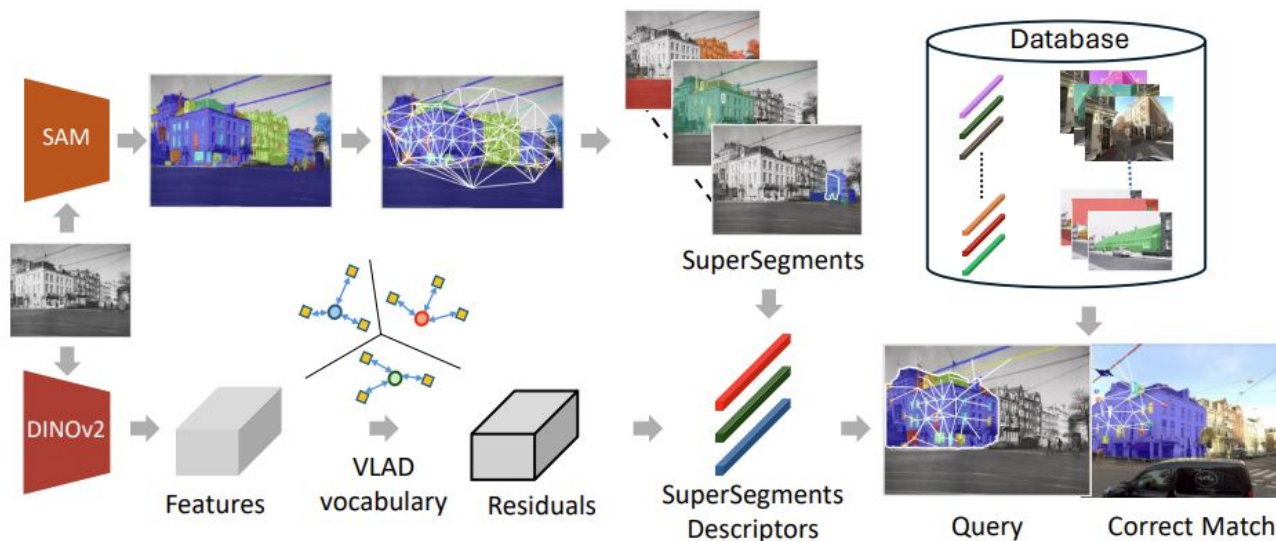
[Paper](#)

[Code](#)

Prior Work From RRC:

Revisit Anything: Visual Place Recognition via Image Segment Retrieval

Kartik Garg^{*1} , Sai Shubodh Puligilla^{*2} , Shishir Kolathaya¹ ,
Madhava Krishna² , and Sourav Garg³ 



[Paper](#)

[Code](#)

Thank You!