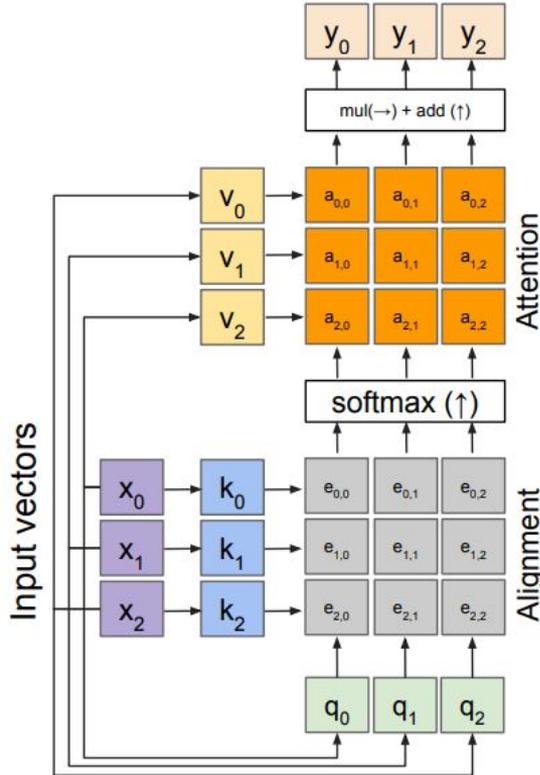# LSTMs and Transformers

CS 231N

# Why Transformers?

Self-Attention is a unique and powerful operation!

# Recap: What is self attention?



**Outputs:**
context vectors: **y** (shape: $D_v$)

**Operations:**
Key vectors: $\mathbf{k} = \mathbf{x}\mathbf{W_k}$
Value vectors: $\mathbf{v} = \mathbf{x}\mathbf{W_v}$
Query vectors: $\mathbf{q} = \mathbf{x}\mathbf{W_q}$
Alignment: $e_{i,j} = q_j \cdot k_i / \sqrt{D}$
Attention: $\mathbf{a} = \text{softmax}(\mathbf{e})$
Output: $y_j = \sum_i a_{i,j} v_i$

**Inputs:**
Input vectors: $\mathbf{x}$ (shape: N x D)

Self-attention allows the network to learn relationships between each segment of the input!

Interesting Application

# See, Hear, and Feel:
# Smart Sensory Fusion for Robotic Manipulation

Hao Li[1]*    Yizhi Zhang[1]*    Junzhe Zhu[1]    Shaoxiong Wang[2]    Michelle A Lee[1]
Huazhe Xu[1]    Edward Adelson[2]    Li Fei-Fei[1]    Ruohan Gao[1]†    Jiajun Wu[1]†
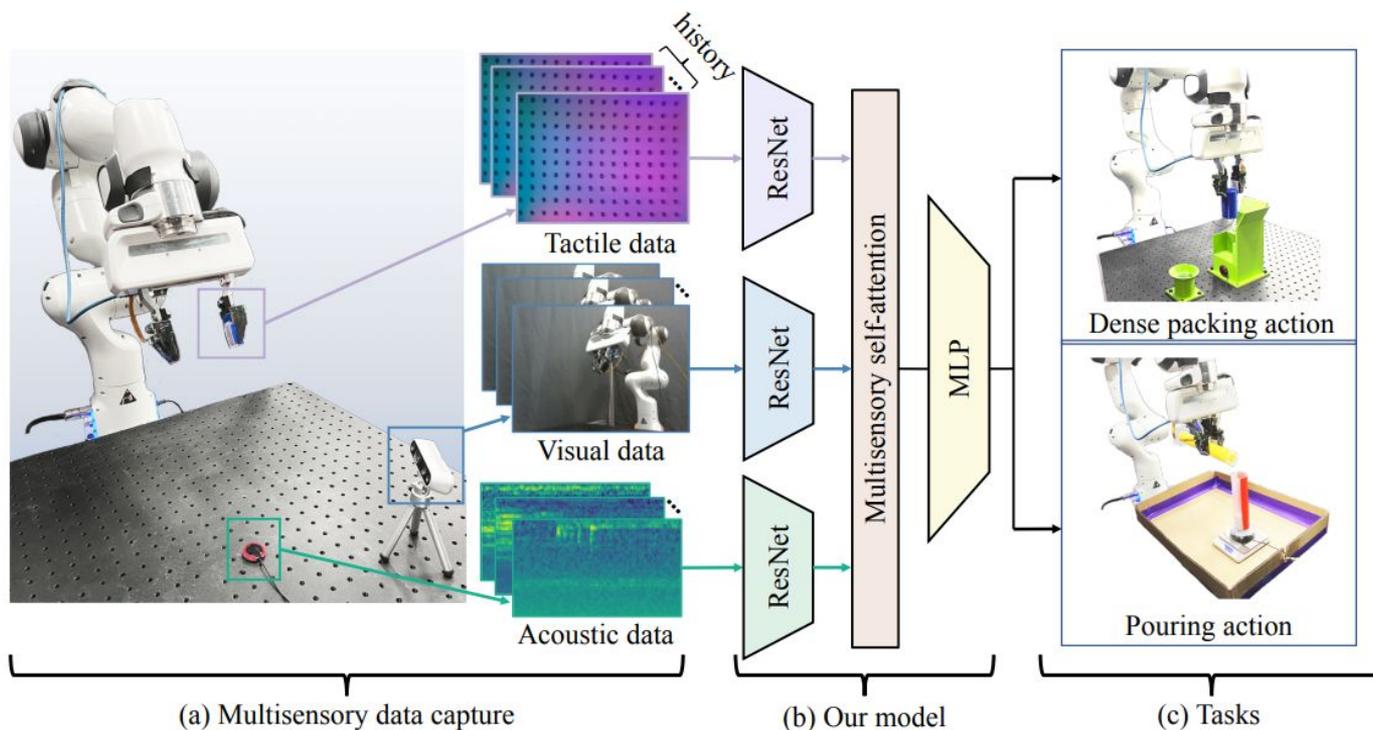[1]Stanford University        [2]Massachusetts Institute of Technology
*Equal contribution.        †Equal advising.
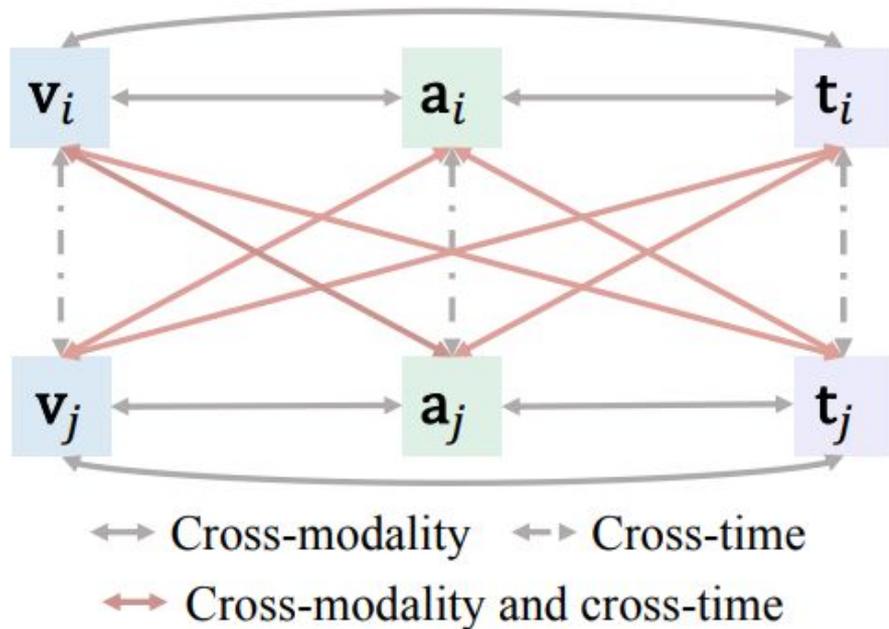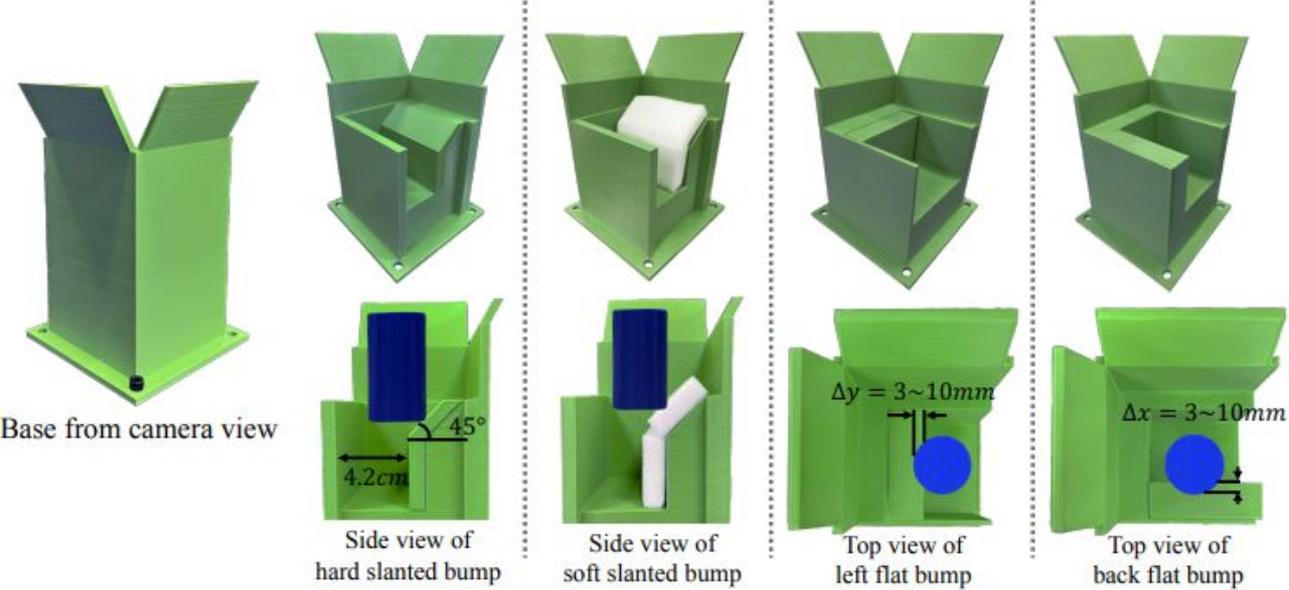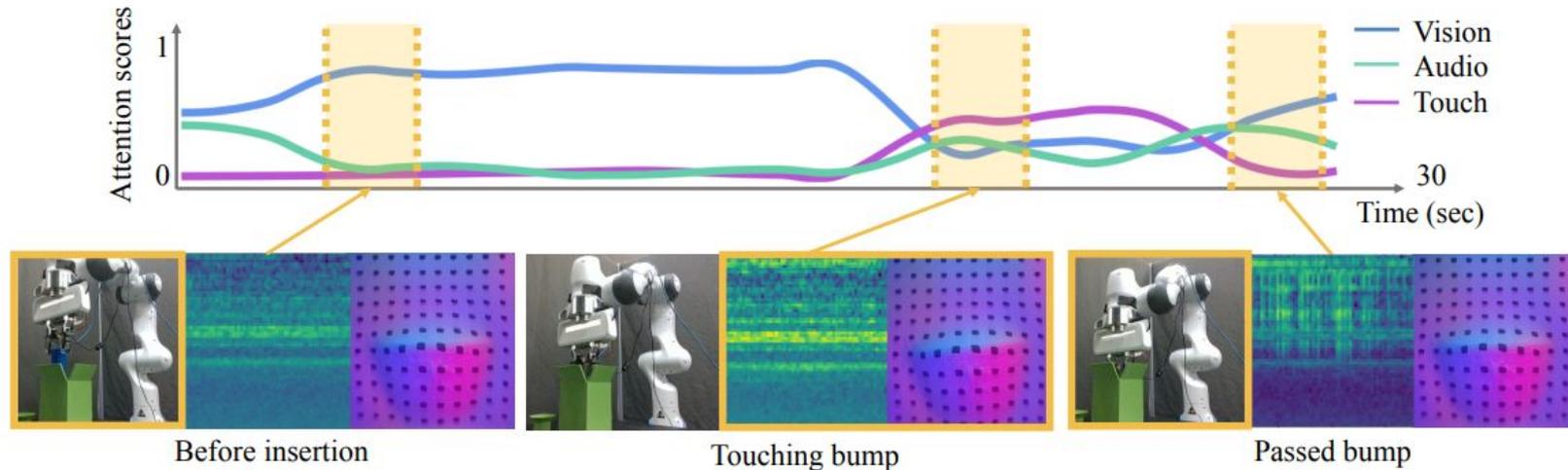
# Overall Setup



(a) Multisensory data capture

Tactile data

Visual data

Acoustic data

history

ResNet

ResNet

ResNet

Multisensory self-attention

MLP

(b) Our model

Dense packing action

Pouring action

(c) Tasks

# Multisensory Self-Attention



Figure 3: Multisensory self-attention.

# Dense Packing



Base from camera view

Side view of hard slanted bump

$4.2 cm$   $45°$

Side view of soft slanted bump

Top view of left flat bump

$\Delta y = 3\sim10mm$

Top view of back flat bump

$\Delta x = 3\sim10mm$

# Tasks Evaluated: Dense Packing



Before insertion

Touching bump

Passed bump

# Real World Example:

CNN vs Transformer

# Examining the Difference Among Transformers and CNNs with Explanation Methods

Mingqi Jiang, Saeed Khorram, Li Fuxin
Collaborative Robotics and Intelligent Systems (CoRIS) Institute
Oregon State University
{jiangmi, khorrams, lif}@oregonstate.edu

# CNN vs Transformer

# Results



(b) ResNet-50-C2

(d) Swin-T