# VIS 2025

# Generalization of CNNs on Relational Reasoning with Bar Charts

Zhenxing Cui*, Lu Chen*, Yunhai Wang, Daniel Haehn, **Yong Wang**, Hanspeter Pfister

# Background

- Convolutional Neural Networks (**CNNs**) are widely used for many visualization tasks
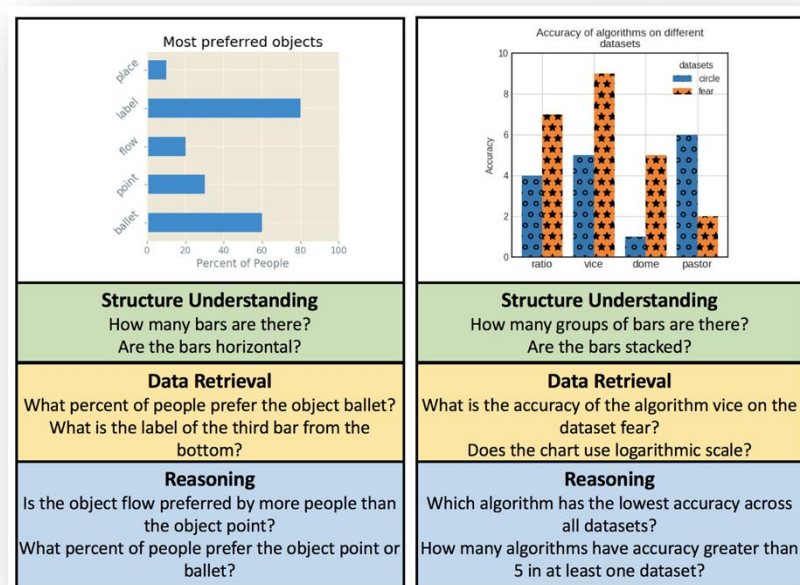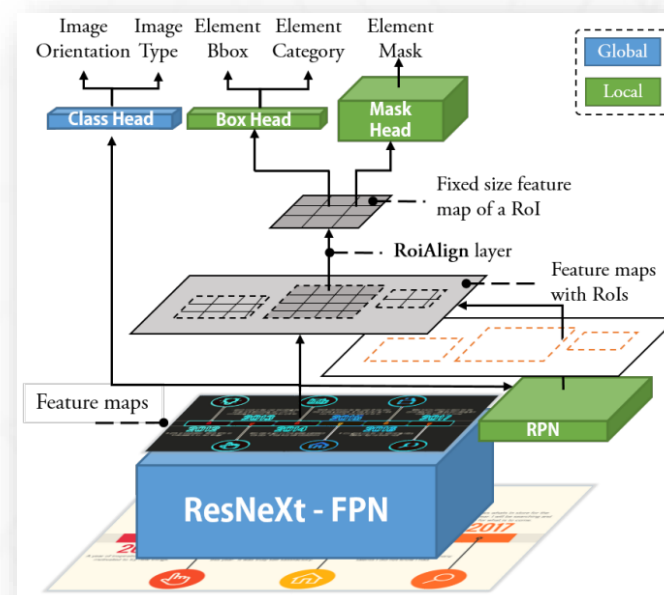


Chart Question Answering



Automated Chart Design
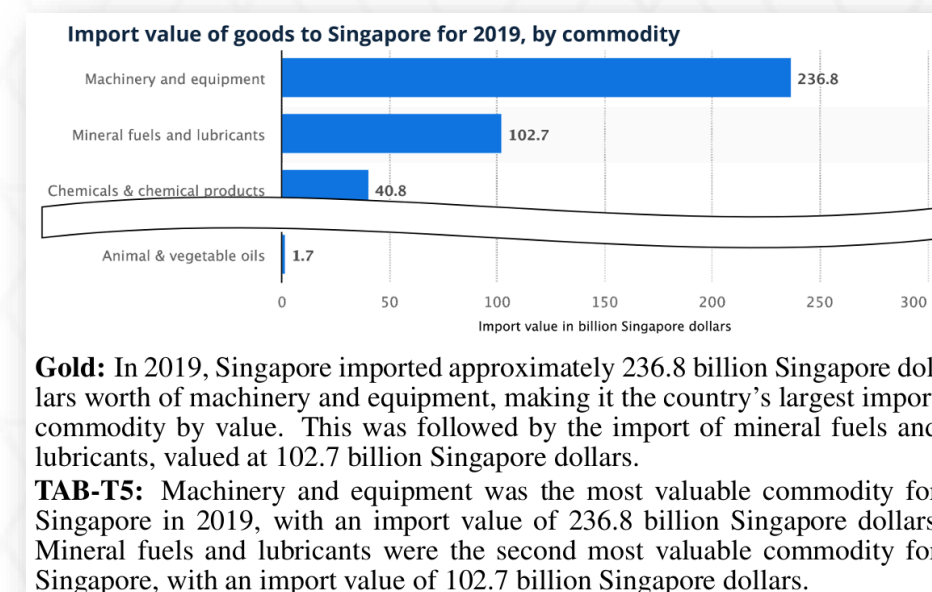


Chart Captioning

- But it remains underexplored how CNNs' graphical perception performance generalizes across visualization design variations

VIS2025

# Background

- **Graphical Perception**: the ability to decode visually encoded quantities in visualizations


Cleveland & McGill, 1984


Haehn *et al.*, 2019

3

VIS2025

# Background

- <span style="color:red">Graphical Perception</span>: the ability to decode visually encoded quantities in visualizations



Oversimplified Charts

Standard Visualizations

VIS 2025

# Research Questions

- How well do CNNs perform on standard visualizations with full design elements?

- How robust are CNNs to design perturbations such as color jitter?

- What are the differences between CNNs and humans in visual relational reasoning?

# Relational Reasoning in Graphical Perception

- **Task**: Estimate the ratio of lengths (i.e., heights) between two target bars (targets indicated by black dots)

# Benchmarking Representative CNNs

- Replicate Haehn *et al.*'s experiments [1] with systematically-tuned CNNs
- CNNs achieve very strong performance, better than previously-reported results

**Architectures**

- MLP
- AlexNet
- LeNet
- DenseNet
- VGG19
- ResNet152
- Xception126
- EfficientNet

**Optimizers**

- AdamW
- SGDM

**Hyper-parameters**

- Learning rate
- Momentum
- Weight decay



[1] D. Haehn, J. Tompkin, and H. Pfister. Evaluating 'graphical perception' with CNNs. IEEE Transactions on Visualization and Computer Graphics, 25(1):641–650, 2019.

# GRAPE: A GRAphical PErception Dataset

- Five bar-chart types synthesized programmatically with Vega-Lite

- 766K (500K for training & 266K for testing) standard visualizations

- Large and controllable dataset to manipulate design parameters



(a) Type 1    (b) Type 2    (c) Type 3    (d) Type 4    (e) Type 5

# Perturbation Setup

- **9 visual parameters**: title position, title size, background color, bar color, stroke color, bar width, stroke width, bar length, dot position

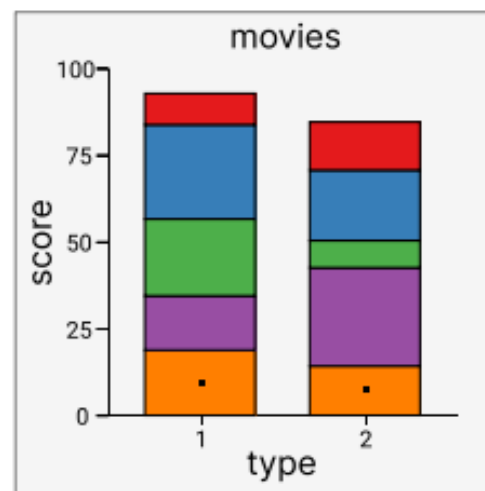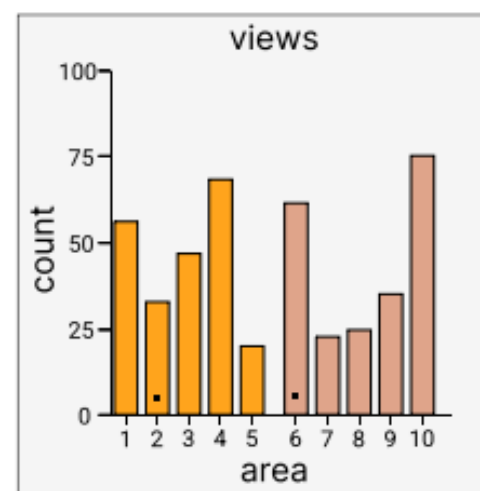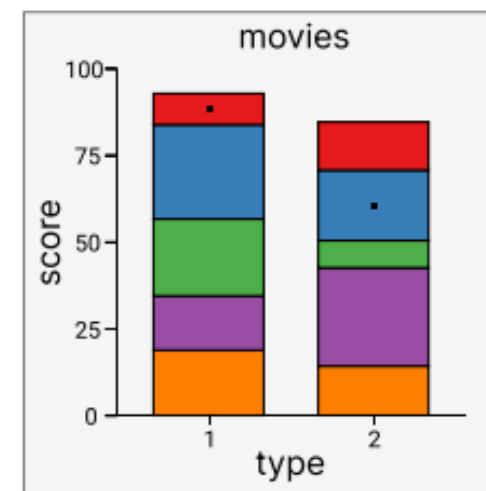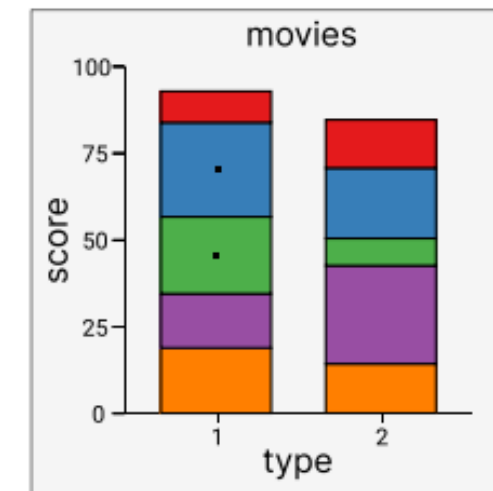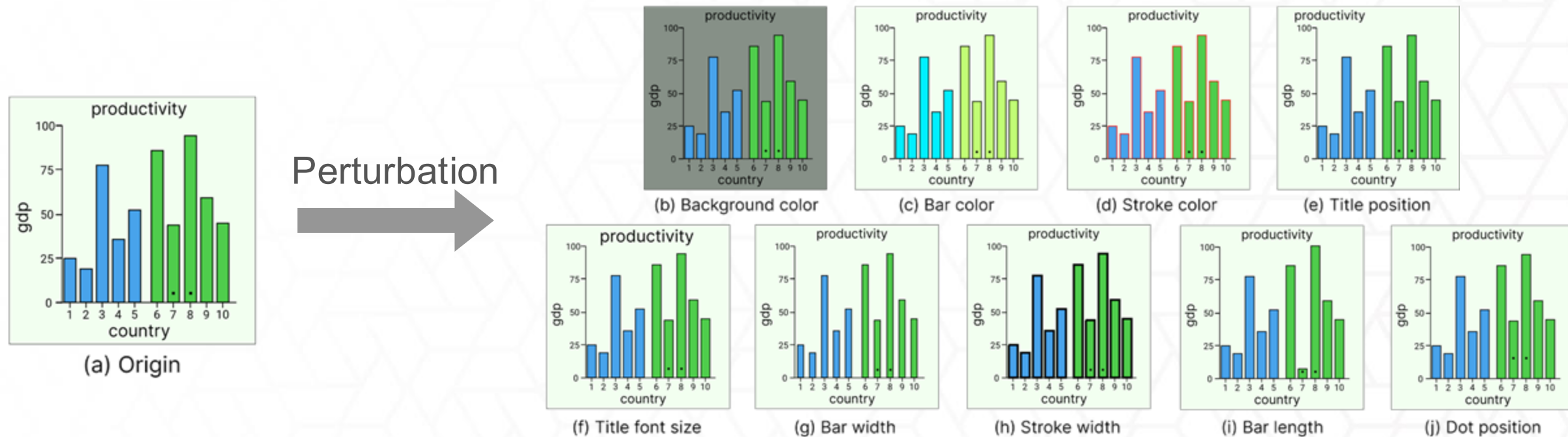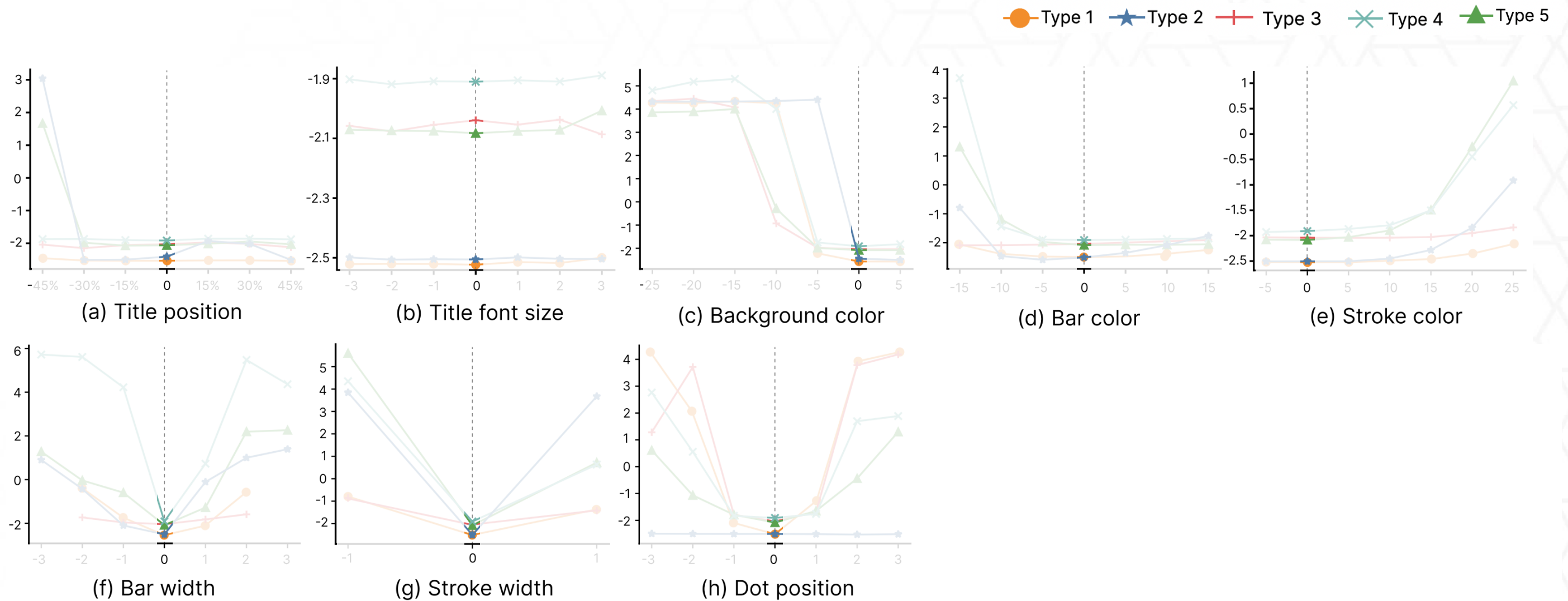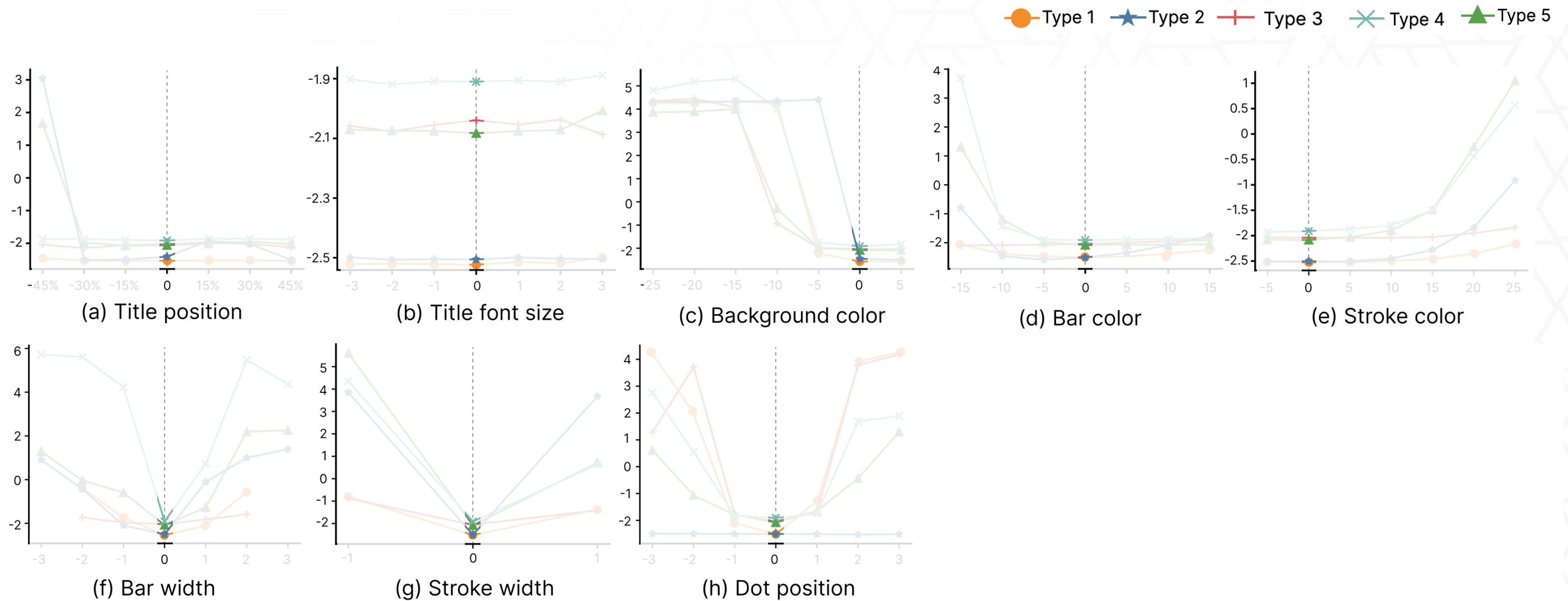- **Independent and Identically Distributed (IID)**: test visualizations have similar encodings with the training samples

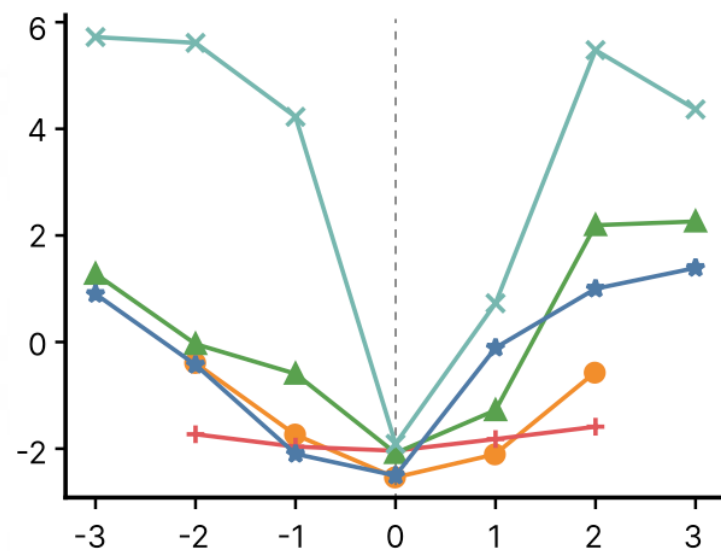- **Out-of-distribution (OOD)**: test and training visualizations are different



(a) Origin

Perturbation

(b) Background color

(c) Bar color

(d) Stroke color

(e) Title position

(f) Title font size

(g) Bar width

(h) Stroke width

(i) Bar length

(j) Dot position

VIS2025

(a) Title position
(b) Title font size
(c) Background color
(d) Bar color
(e) Stroke color
(f) Bar width
(g) Stroke width
(h) Dot position

Type 1  Type 2  Type 3  Type 4  Type 5

VIS 2025

# Results (OOD):
# CNN Robustness Collapses on Visual Parameter Shifts



(a) Title position  (b) Title font size  (c) Background color  (d) Bar color  (e) Stroke color

(f) Bar width  (g) Stroke width  (h) Dot position

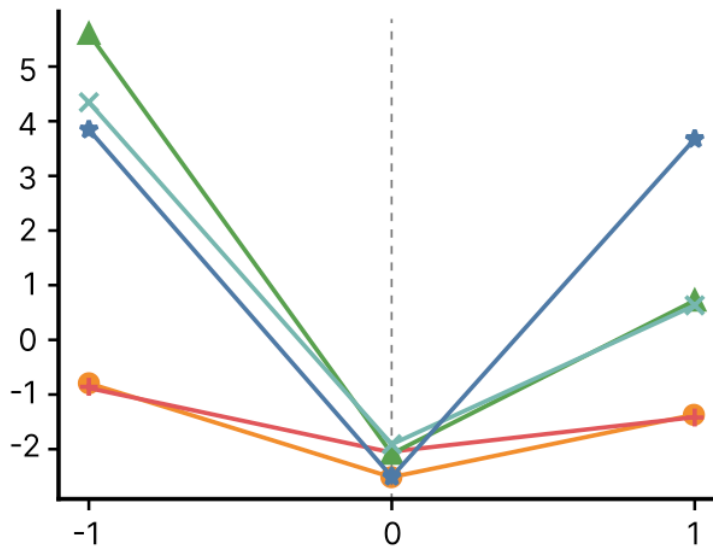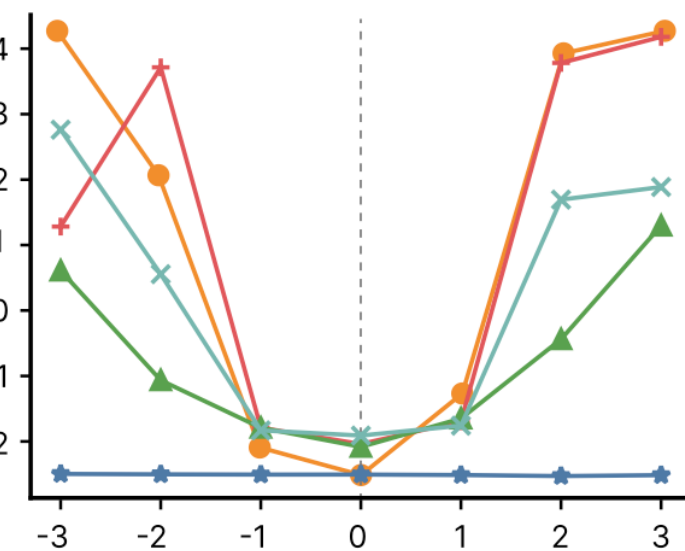# Results (OOD):
# CNN Robustness Collapses on Visual Parameter Shifts

- Small changes in **bar width, stroke width**, and **dot position** → Big errors



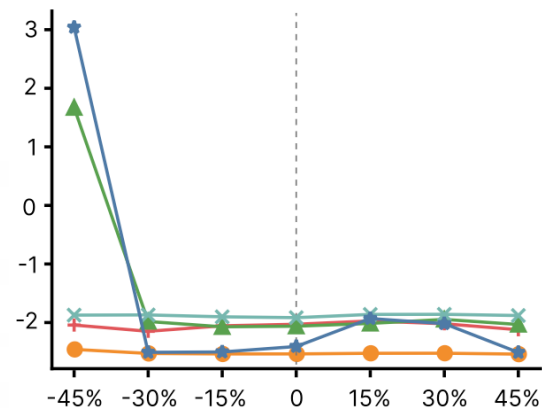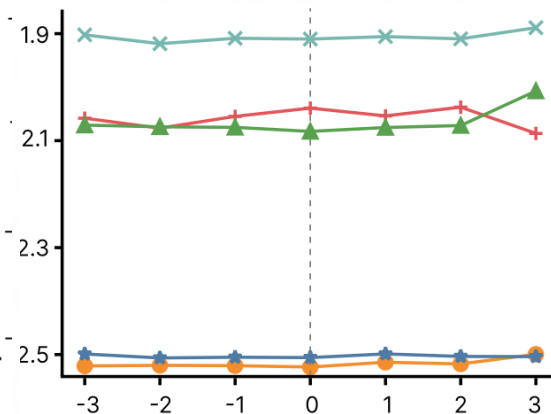(f) Bar width  (g) Stroke width  (h) Dot position

VIS 2025

# Results (OOD):
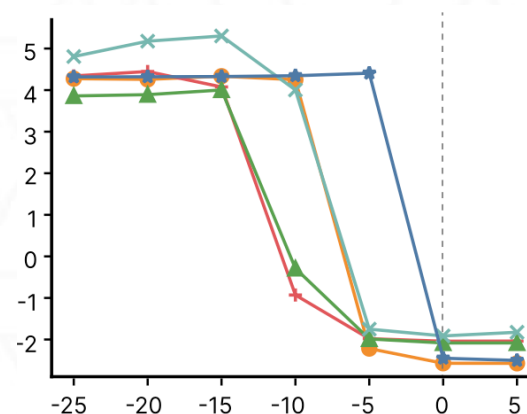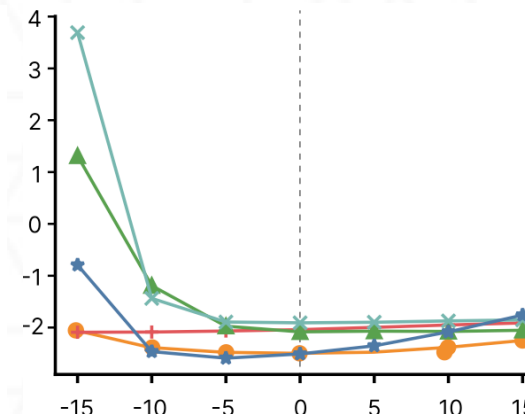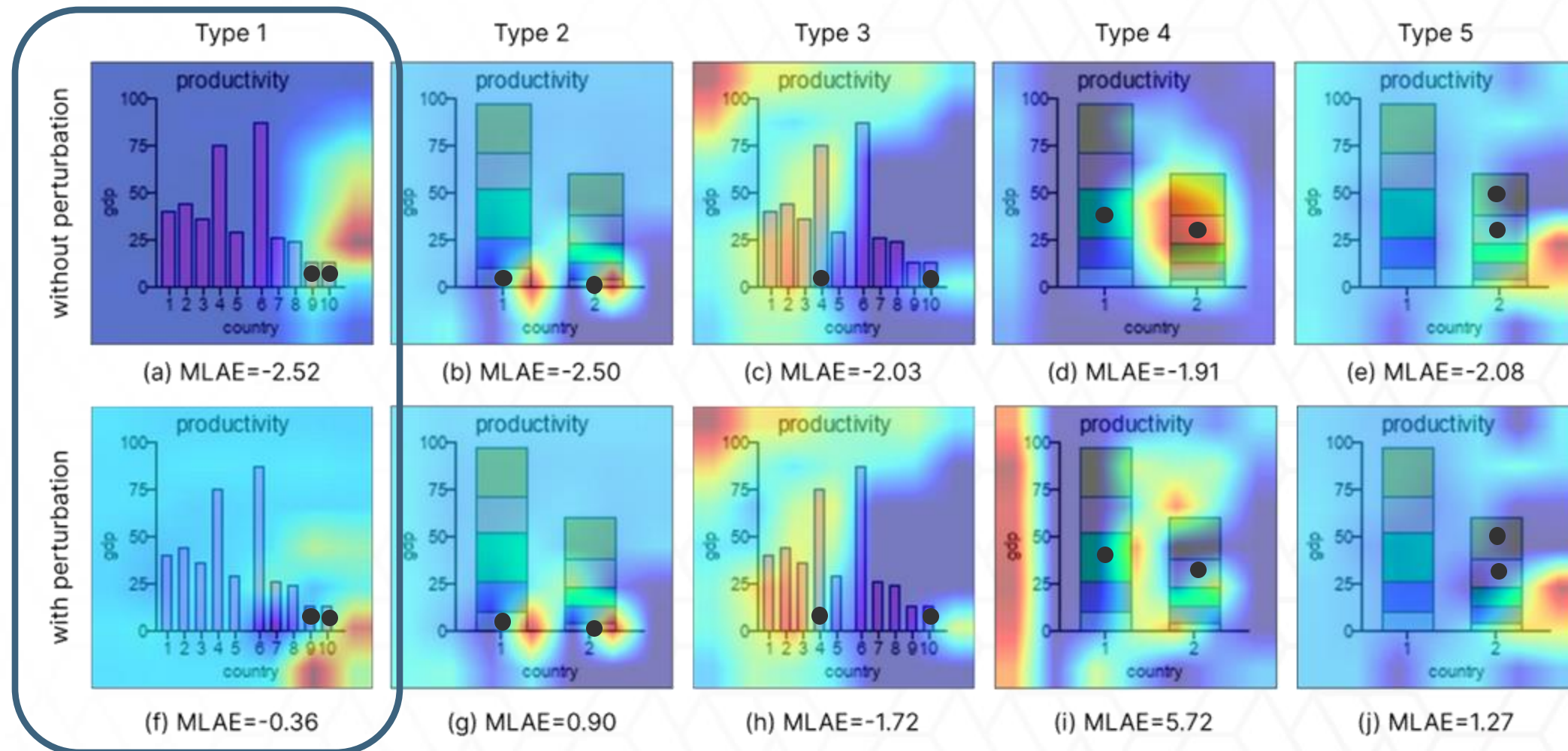# CNN Robustness Collapses on Visual Parameter Shifts

- Even "irrelevant" parameters (e.g., **title position**) affect graphic perception performance

- **Background luminance** drops sharply push up prediction error

- CNNs are relatively robust to changes of **title font size**, **bar and stroke colors**



(a) Title position    (b) Title font size    (c) Background color    (d) Bar color    (e) Stroke color

# Humans vs CNNs

- Humans are worse than CNNs in IID conditions but more robust under perturbations (OOD)

- Interview feedback: participants focus on target bars; ignore nonessential styling

# Why Do CNNs Fail OOD? Attention on the Wrong Pixels

- Grad-CAM saliency maps rarely highlight target bars
- Model attention shifts noticeably under minor perturbations of bar width

# Can We Fix It? Segmentation Masks

- Add a target-bar mask channel (RGB-α) → Better target localization, some robustness gains

- Yet, still sensitive to shape-related visual parameters like bar and stroke width

# Does Data Augmentation Solve It?

- Performance improves on those perturbations seen at train time
- Generalization to unseen perturbations remains weak

# Take-Home Messages

- Small and unseen shifts of visual parameters break CNNs' graphical perception performance

- CNNs still fail to see charts like humans

- Simple target masks and dataset augmentation aren't enough for enhancing the generalization of CNNs

VIS2025

# Call for More Research on Benchmarking AI4Vis!



Latest AI Models

Different Visualization Tasks

# Generalization of CNNs on Relational Reasoning with Bar Charts

Zhenxing Cui*, Lu Chen*, Yunhai Wang, Daniel Haehn, **Yong Wang**, Hanspeter Pfister