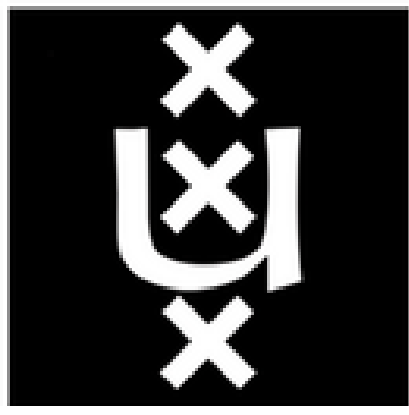


# Using Explainable AI to Investigate Navigational Affordances in Real-world Scenes



UNIVERSITY  
OF AMSTERDAM

Clemens G. Bartnik and Iris I.A. Groen

Video & Image Sense Lab, Informatics Institute, University of Amsterdam, The Netherlands

**Aims** We move with ease through our environment. We can choose different paths and use a variety of different navigational actions, such as walking, swimming or climbing. Past research<sup>1</sup> focusing on indoor environments<sup>2</sup> suggests that CNNs<sup>3</sup> trained for scene recognition<sup>4</sup> contain useful representations of navigational affordances.

Here, we investigated the ability of CNNs to capture affordances in a broader set of environments and use different explainable AI feature visualizations (LRP<sup>5</sup>, LIME<sup>6</sup>, Grad-CAM<sup>7</sup>) and different task objectives/network models (depth perception<sup>8</sup>, scene segmentation<sup>9</sup>).

## Method

### XAI feature visualizations:

**Layer-wise Relevance Propagation (LRP)**<sup>5</sup> identifies important pixels – those that contribute the most to the prediction get a higher relevance score.

**Local Interpretable Model-Agnostic Explanations (LIME)**<sup>6</sup> creates local explanations by perturbing the input and tracking how the predictions change.

**Grad-CAM**<sup>7</sup> visualizes image regions that were important for the classification decision.

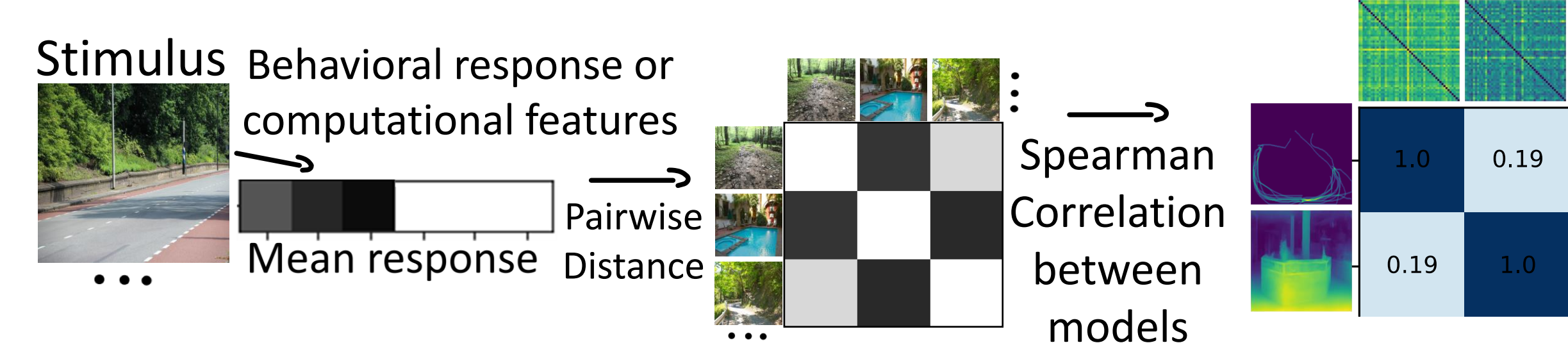
### Models:

**VGG16**<sup>3</sup> trained for scene recognition on Places 365<sup>4</sup>

**Scene Segmentation**<sup>9</sup> trained on ADE20K – for stimuli<sup>2</sup> we use floor segmentation. For our stimuli we selected the segmentation class with the highest pixel overlap with the average paths.

**Depth estimation**<sup>8</sup> resulting depth map of monocular depth estimation model

### Representational similarity analysis (RSA) <sup>11</sup>



Distance metric: Euclidean distance

## Behavior

In an online experiment we collected human annotations (N = 152) of possible paths through each scene (at least 20 per image) like <sup>2</sup>. Each participant was allowed to draw only one possible path but was not otherwise constrained.

### Task Description

“Draw a route along which you would move/navigate in any way (walk, swim, ride, or climb) through each scene.”

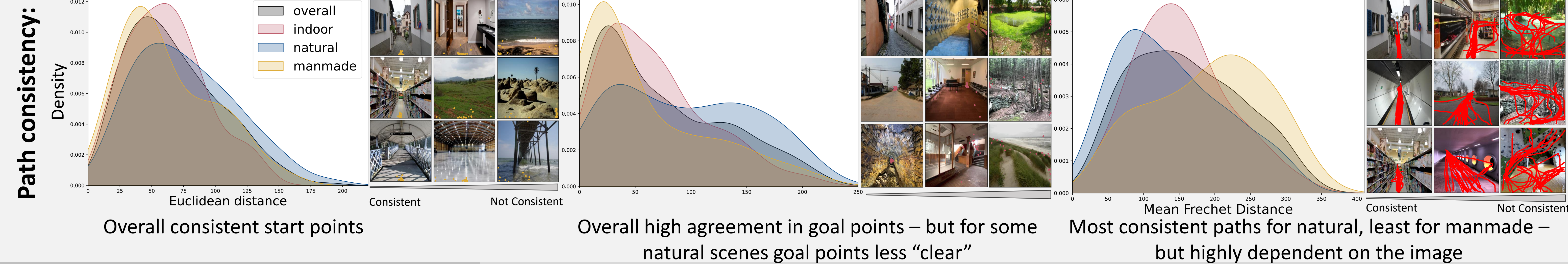
### Trial Example



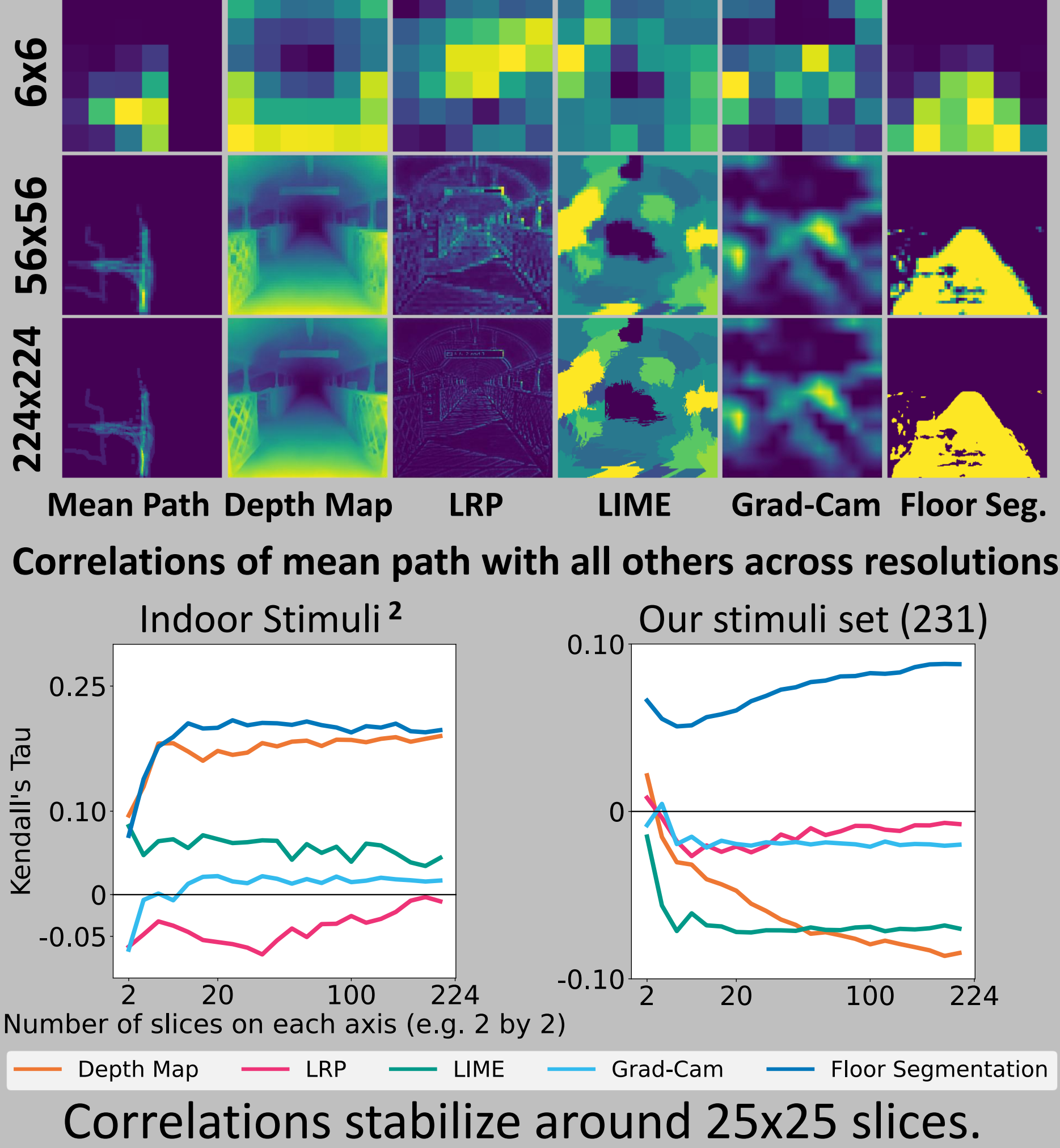
### Average path across participants



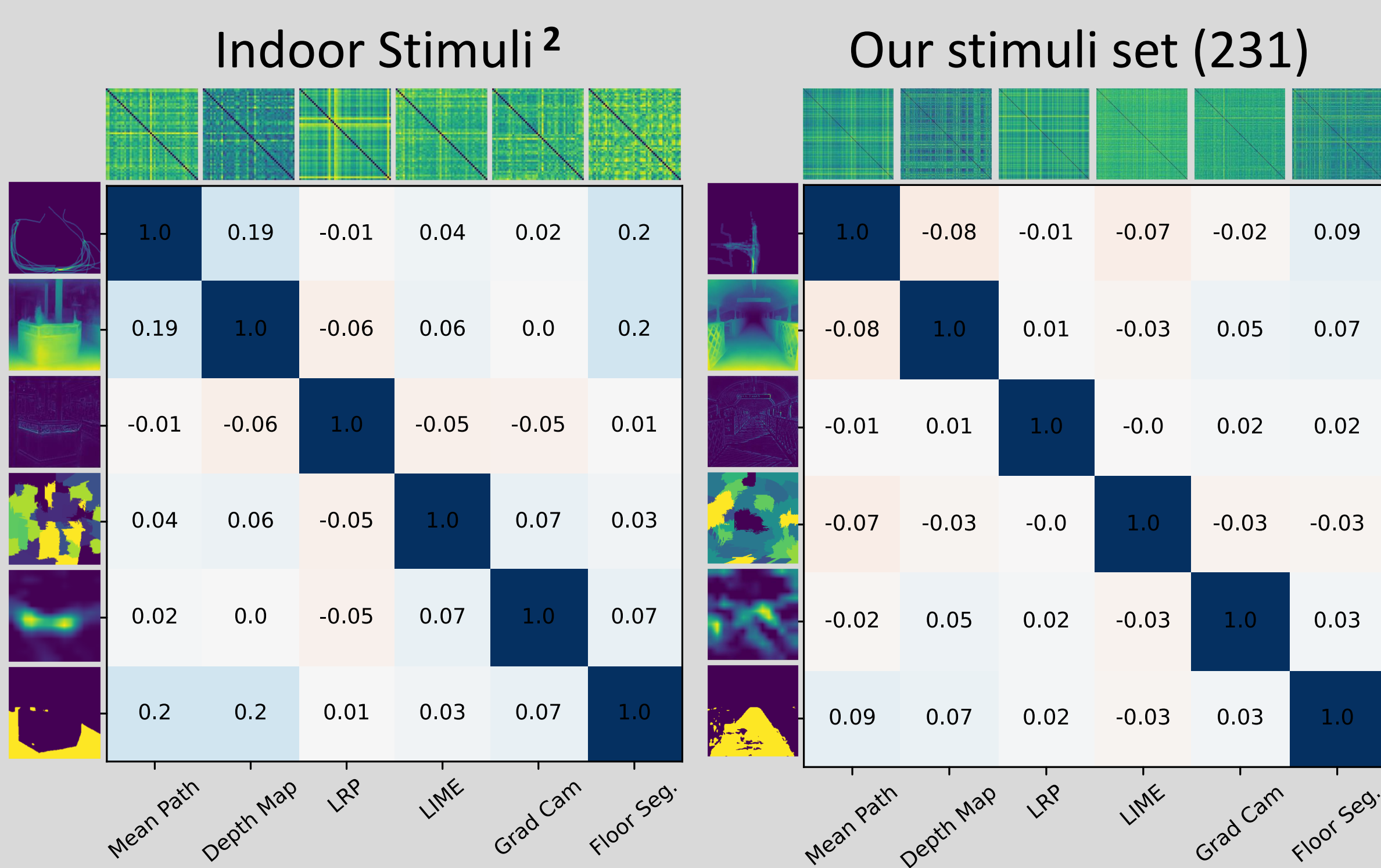
## Results



### Resolution comparisons:

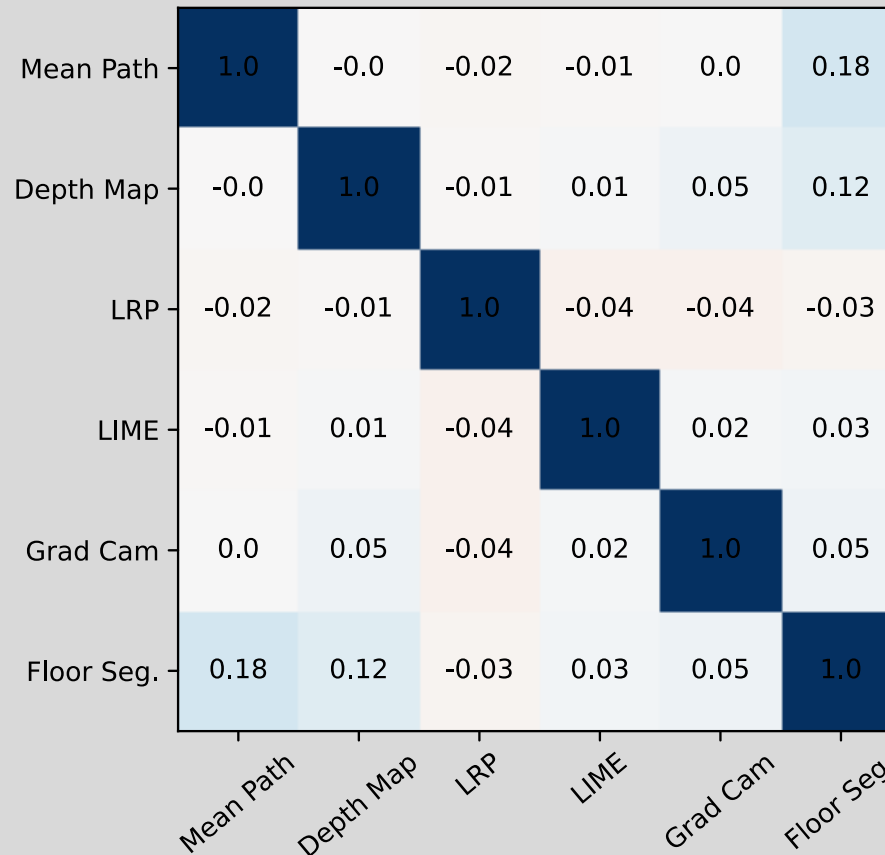


### Model comparisons:

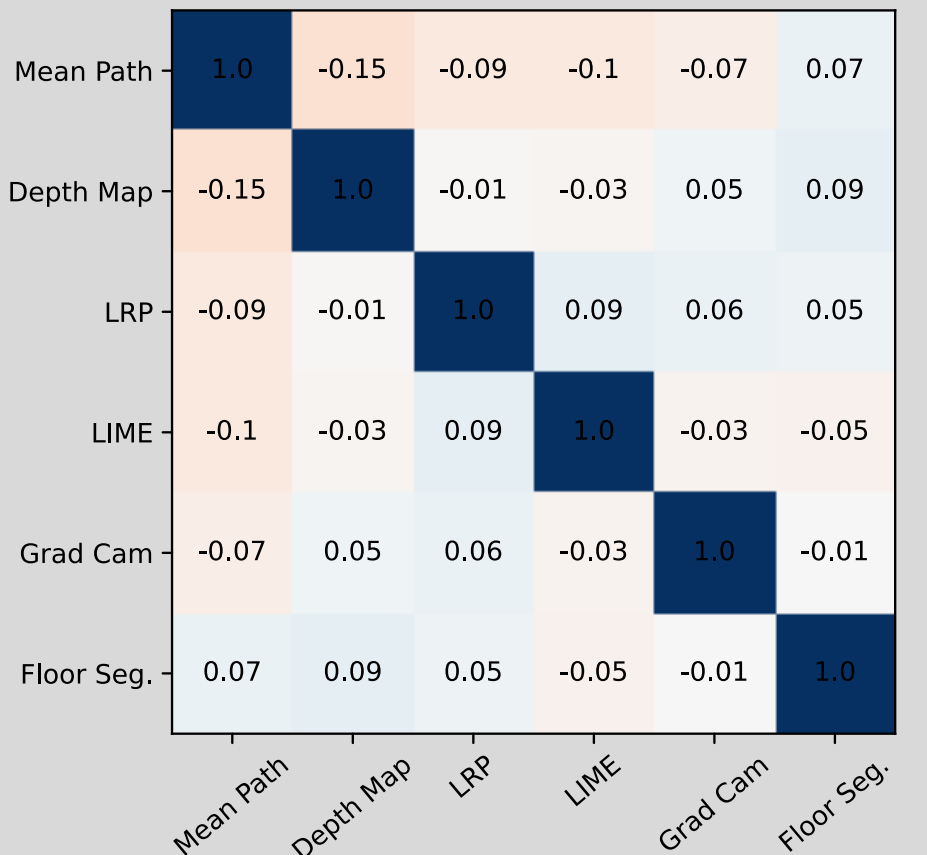


In the stimuli<sup>2</sup> we see substantial correlations between certain feature visualizations. These are much lower in our new stimulus set with diverse environments.

### Indoor

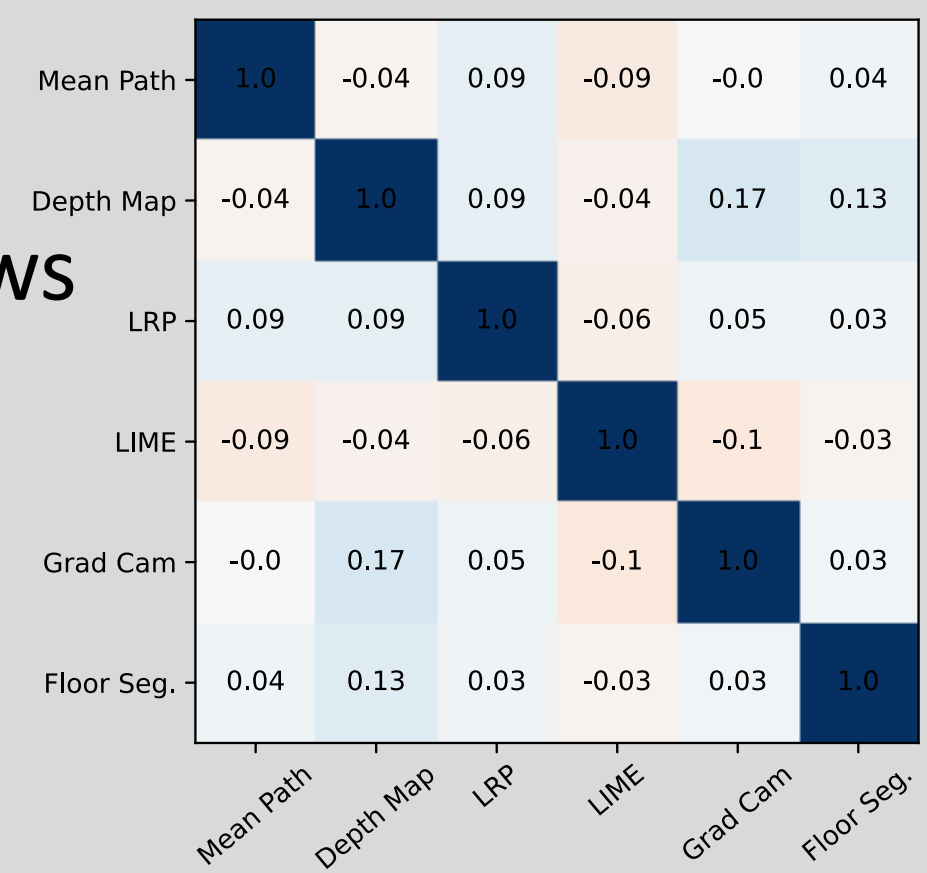


### Outdoor manmade



The correlation pattern changes for different types of environments. Our indoor subset shows similar pattern to stimuli<sup>2</sup> but the other types of environments differ.

### Outdoor natural



## Conclusions

- Human participants are consistent in annotating possible paths in diverse natural scenes. In most scenes they independently pick consistent goal points and use consistent paths to navigate the presented scene. This shows that diagnostic features for potential paths are present.
- We can confirm that scene layout<sup>1</sup> captured by depth estimation and floor segmentations captures affordances of indoor environments. But this does not generalize to other types of environments (outdoor manmade, outdoor natural)..
- None of the CNN feature representation we tested, alone adequately captures navigational affordances across all types of environments. It is still unclear how humans incorporate multiple relevant features to perceive a possible path through a natural scenes.

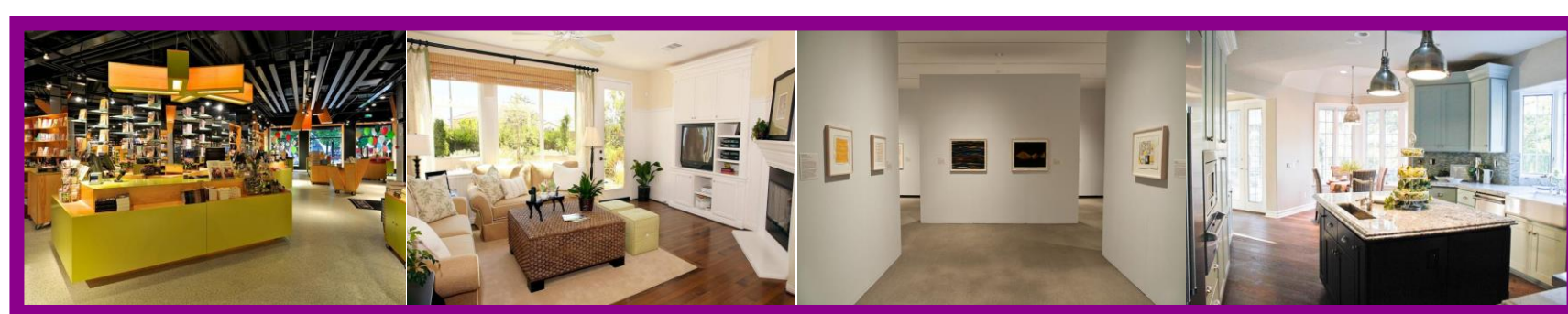
### References

<sup>1</sup> Bonner & Epstein (2018), *PLoS Comput. Biol.* | <sup>4</sup> Zhou et al. (2017), *IEEE PAMI* | <sup>7</sup> Selvaraju et al. (2020) *Int. J. Comput. Vis.* | <sup>10</sup> Patterson & Hays (2012), *CVPR* | <sup>2</sup> Bonner & Epstein (2017), *PNAS* | <sup>5</sup> Bach et al. (2015), *PLoS One* | <sup>8</sup> Miangoleh et al. (2021), *CVPR* | <sup>11</sup> Kriegeskorte et al. (2008), *Front. Syst. Neurosci* | <sup>3</sup> Kalliatakis (2017), *GitHub* | <sup>6</sup> Ribeiro et al. (2016), *Proc. ACM SIGKDD* | <sup>9</sup> Zhou et al. (2018), *Int. J. Comput. Vis.* |

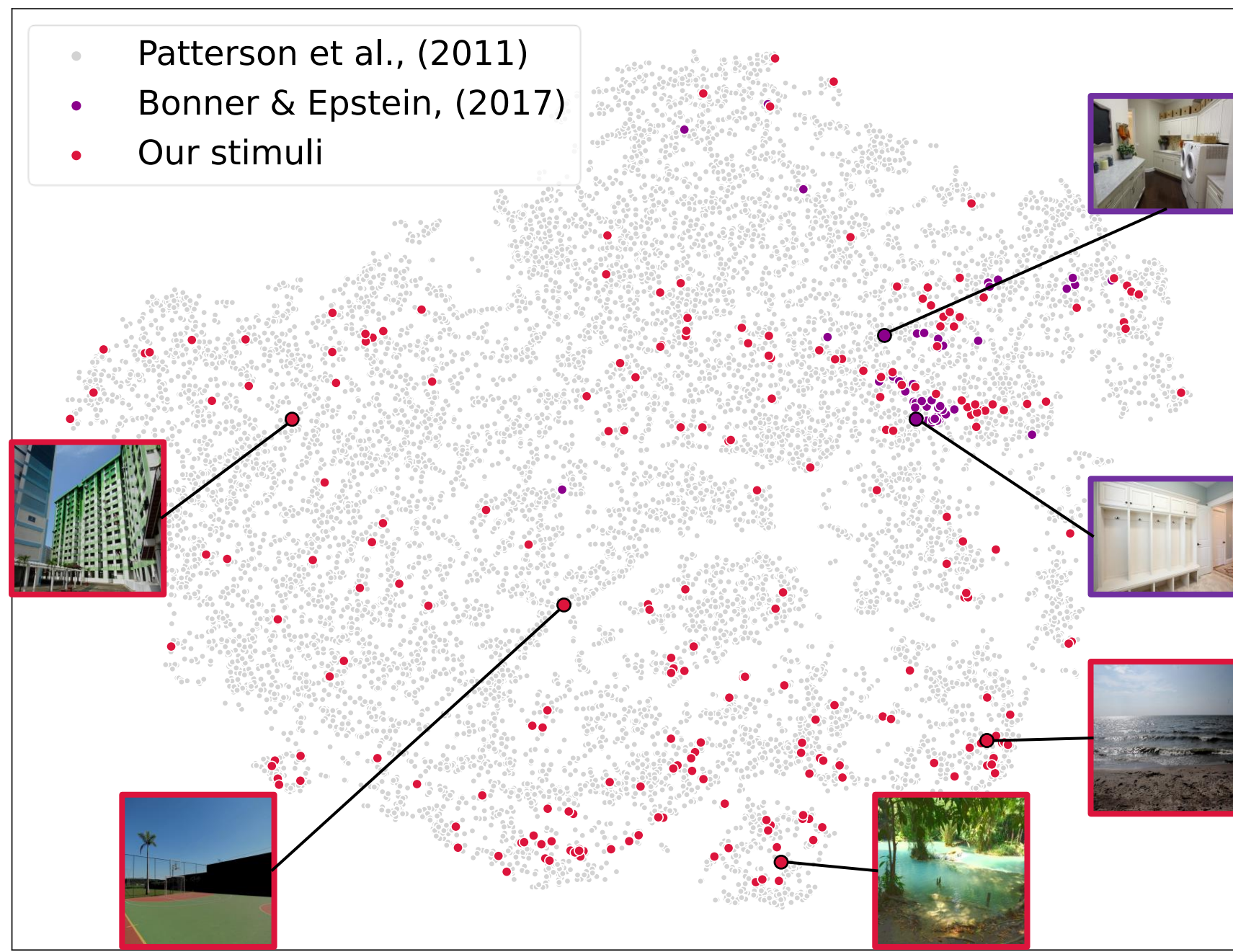
231

## Stimuli

50 naturalistic real-world indoor stimuli<sup>1</sup>



naturalistic real-world stimuli (equal No. of images: indoor, outdoor natural and outdoor manmade)



231