

# Comparing supervised learning dynamics: Deep neural networks match human data efficiency but show a generalisation lag

Lukas S. Huber<sup>1,2</sup>, Fred W. Mast<sup>2</sup>, Felix A. Wichmann<sup>1</sup>

Most  
human-to-DNN  
behavioral comparison  
studies of object  
recognition focus on  
final outcomes, not on  
the learning  
process.

Empirical  
side-by-side comparison  
of supervised  
representation learning in  
humans and DNNs with  
aligned learning  
conditions.

Results  
indicate that while  
DNNs match humans  
in data efficiency,  
they show a  
pronounced  
generalization lag,  
unlike humans.

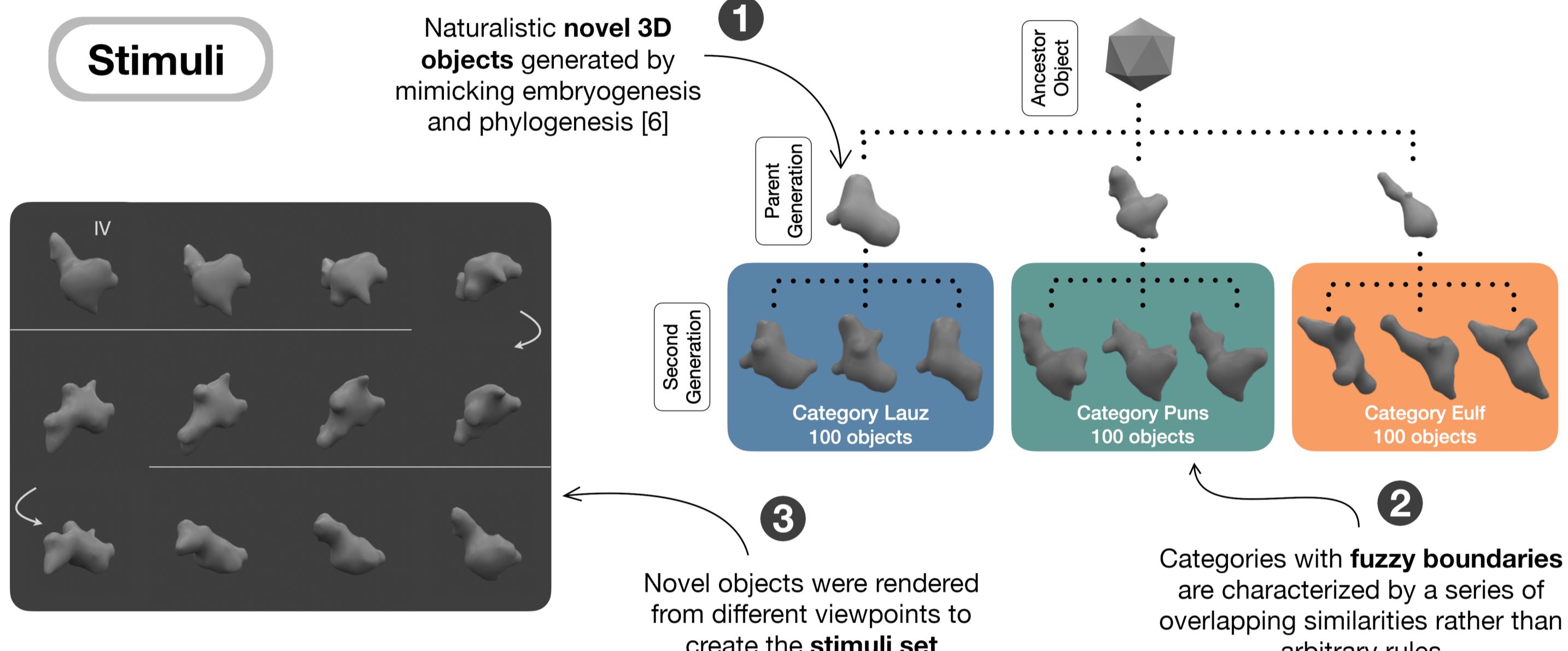
## Motivation

Despite deep neural networks (DNNs) surpassing human-level performance in object recognition tasks [1], more fine-grained behavioral assessments reveal significant representational divergences, as evidenced by differences in robustness [2], error patterns at the category and trial level [3], susceptibility to adversarial attacks [4], and crowding [5]. Identifying the origins of these representational differences is a non-trivial task, given the vast search space for potential causes.

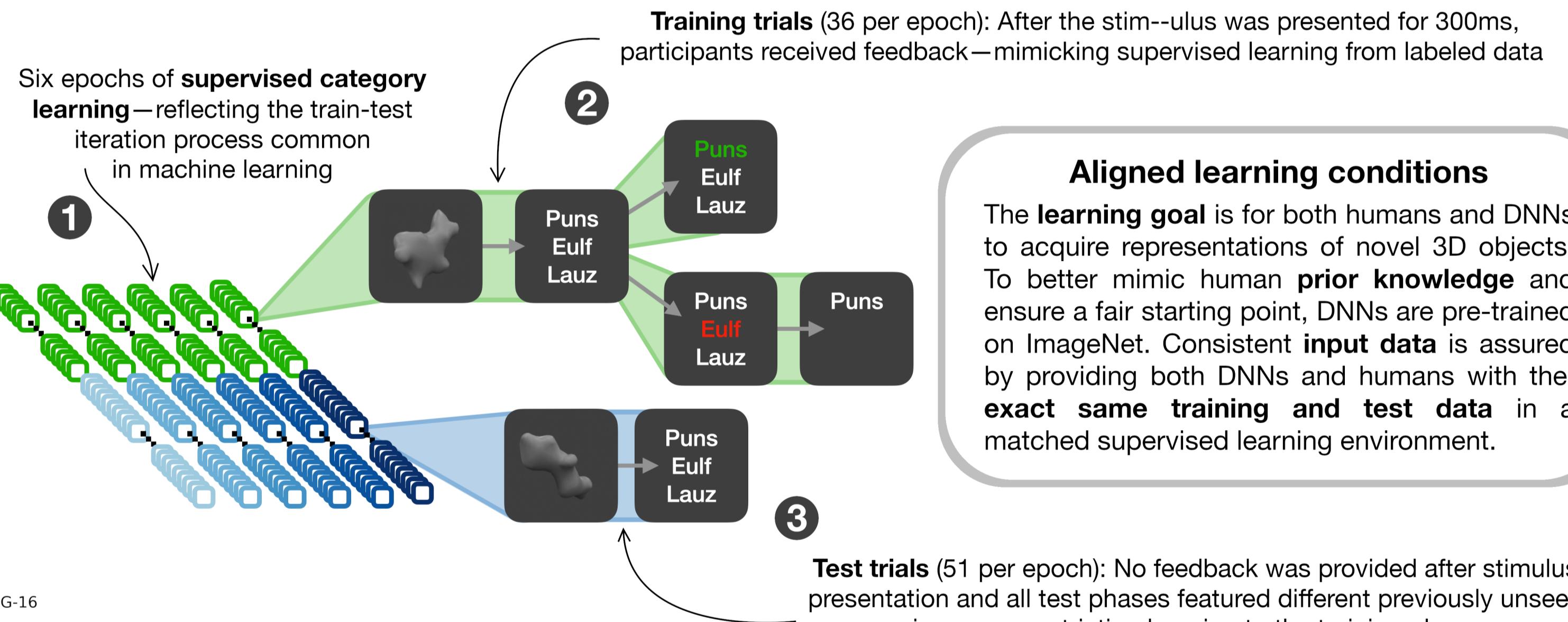
## Rationale

Comparing the process in which representations are acquired appears crucial: for each learning condition aligned during representation acquisition, the search space becomes more constrained. We investigate how humans and DNNs acquire novel representations in a constrained environment with aligned learning conditions, tracking behavioral changes across six epochs of a category learning task. At each epoch, we assess the generalization of learned representations to previously unseen data.

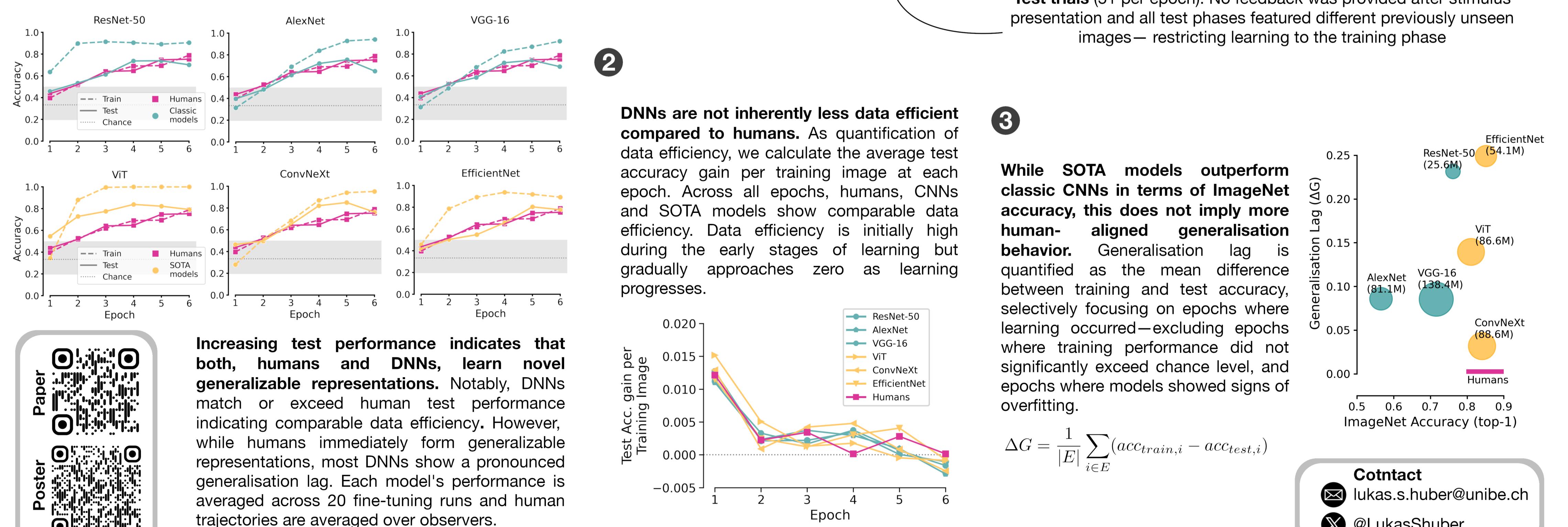
## Stimuli



## Learning environment



## Results



Increasing test performance indicates that both, humans and DNNs, learn novel generalizable representations. Notably, DNNs match or exceed human test performance indicating comparable data efficiency. However, while humans immediately form generalizable representations, most DNNs show a pronounced generalisation lag. Each model's performance is averaged across 20 fine-tuning runs and human trajectories are averaged over observers.

[1] He, K., Zhang, X., Ren, S., & Sun, J. (2015). Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In Proceedings of the IEEE International Conference on Computer Vision.

[2] Geirhos, R., Temme, C. R., Rauber, J., Schütt, H. H., Bethge, M., & Wichmann, F. A. (2020). Beyond accuracy: quantifying trial-by-trial behaviour of CNNs and humans by measuring error consistency. Advances in Neural Information Processing Systems, 33.

[3] Zhou, Z., & Firestone, C. (2019). Humans can decipher adversarial images. *Nature Communications*, 10(1), 1334.

[4] Lonnqvist, B., Bornet, A., Doering, A., & Herzog, M. H. (2021). A comparative biology approach to DNN modeling of vision: A focus on differences, not similarities. *Journal of Vision*, 21(10), 17–17.

[5] Hauffen, K., Bart, E., Brady, M., Kersten, D., & Hegde, J. (2012). Creating objects and object categories for studying perception and perceptual learning. *JoVE (Journal of Visualized Experiments)*, 69, e3538.

