# $\partial$HT
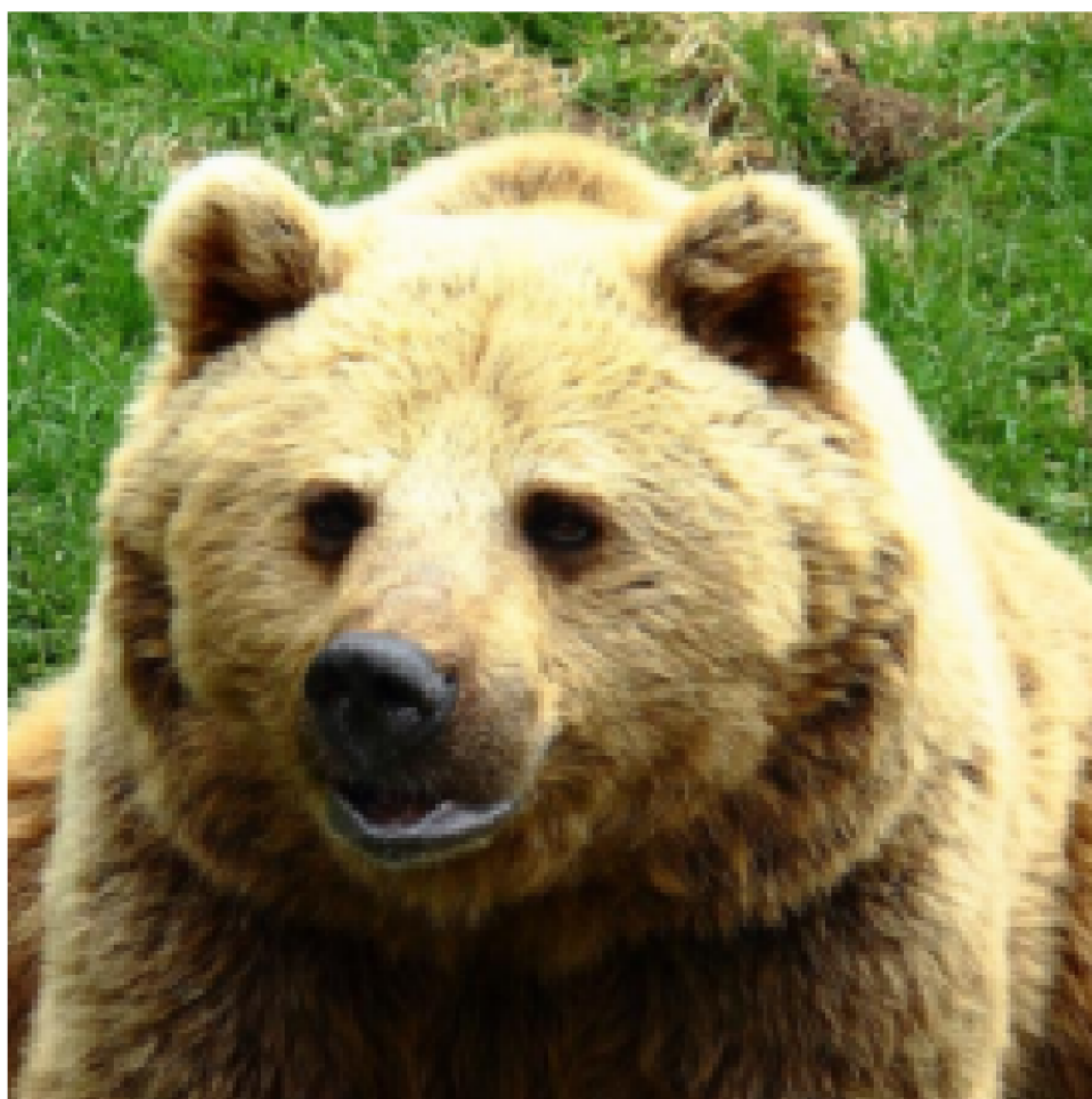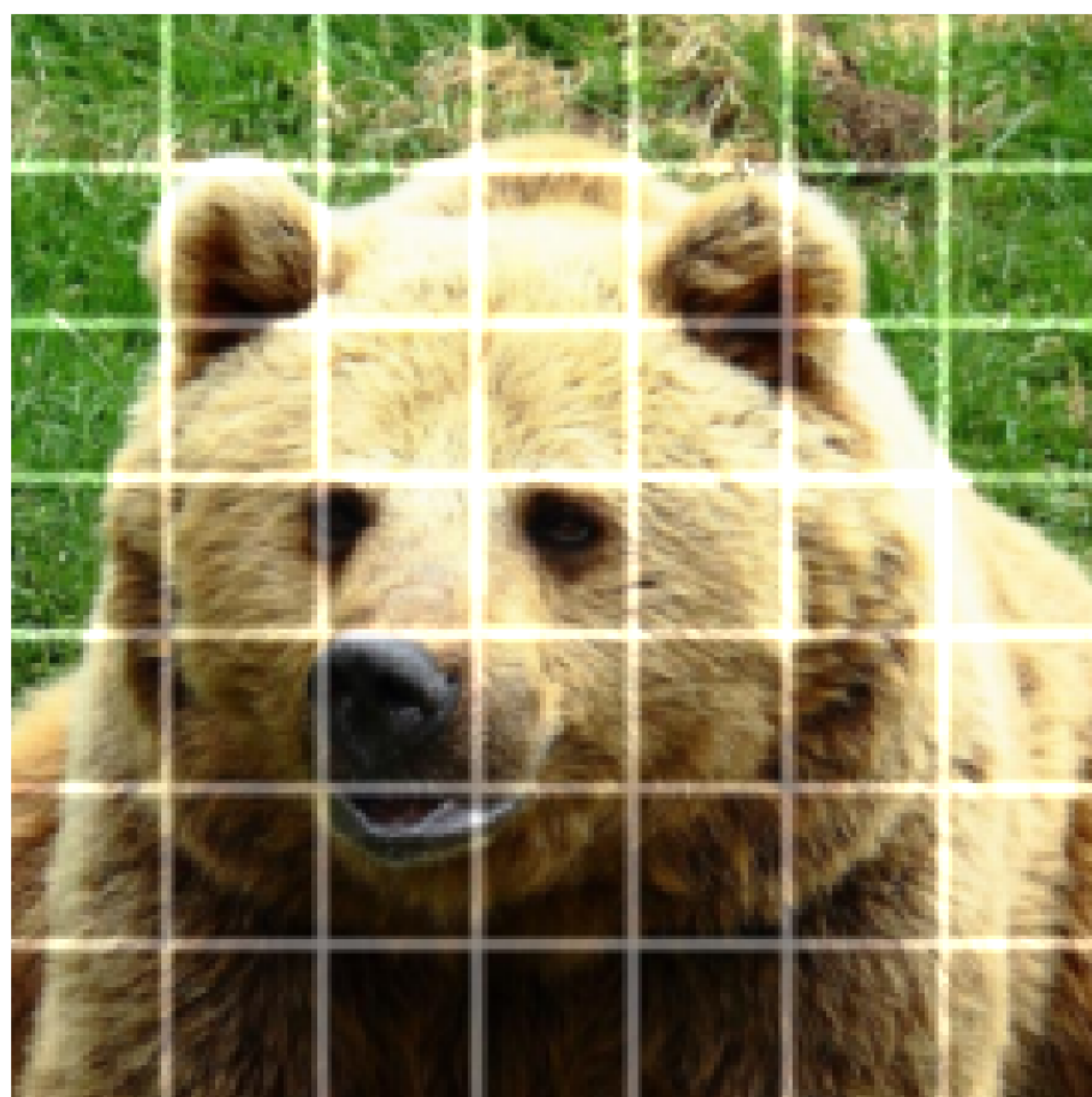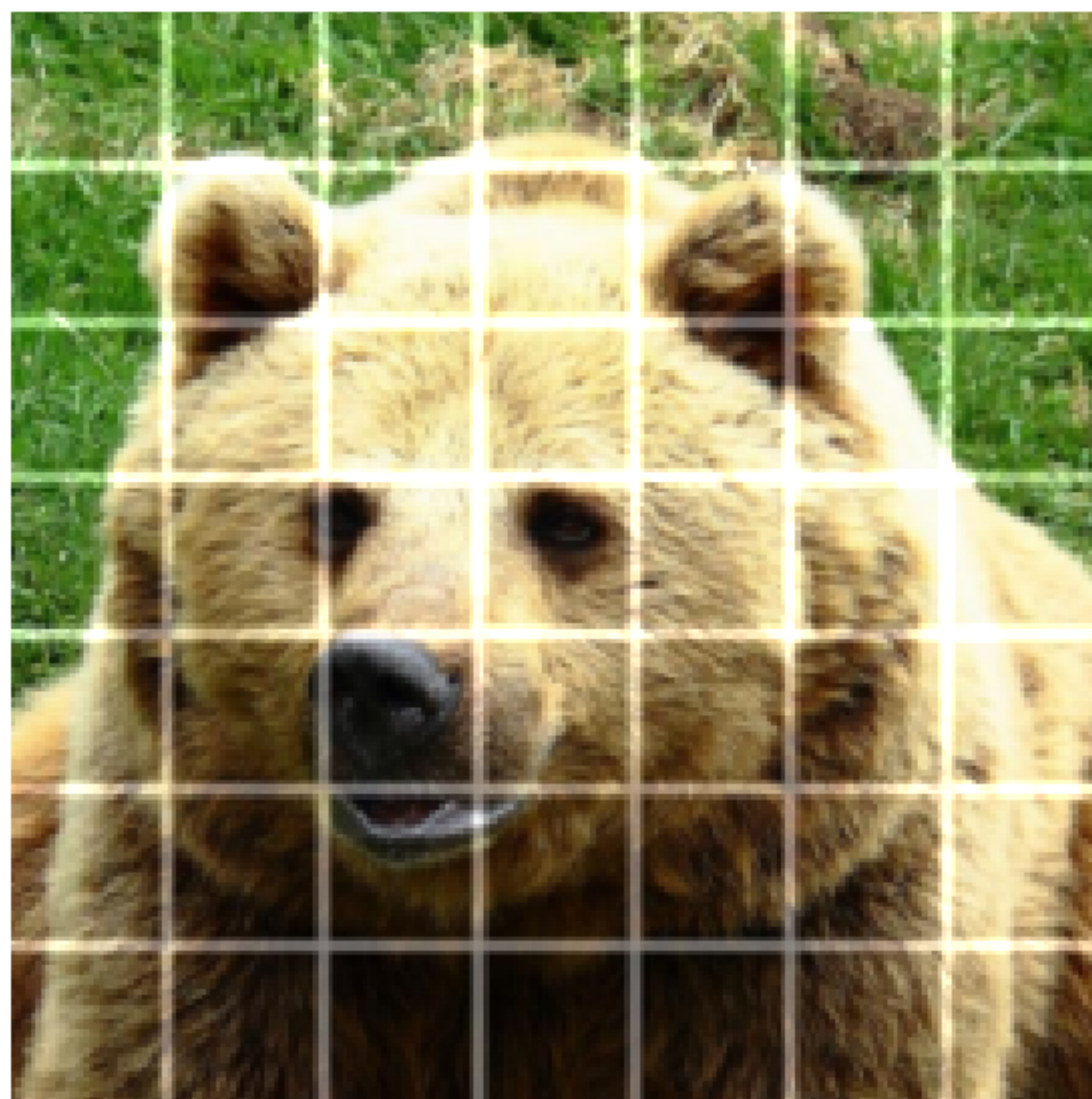## Differentiable Hierarchical Visual Tokenisation

**Marius Aasan - University of Oslo | SFI Visual Intelligence**

UNIVERSITY OF OSLO
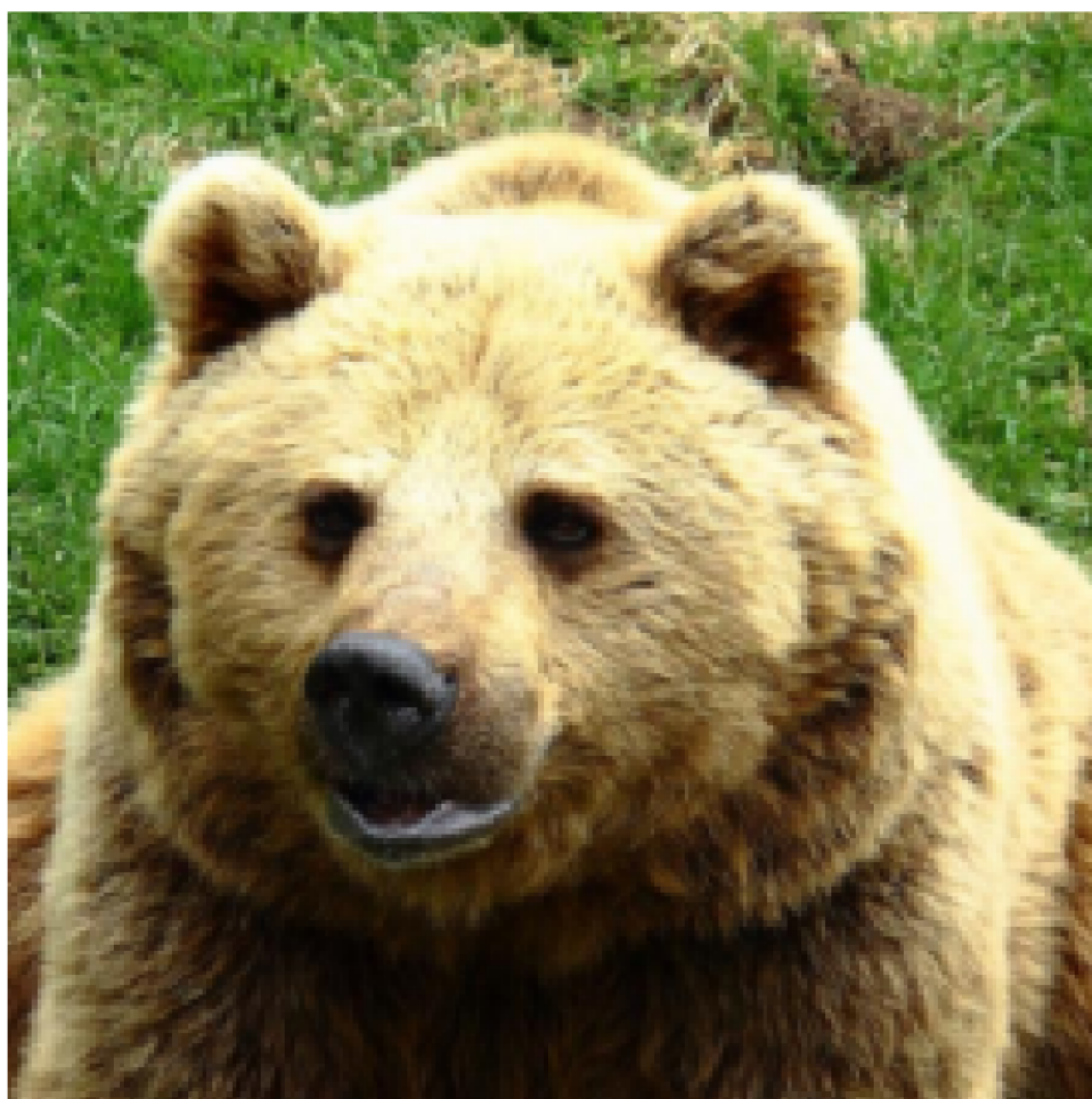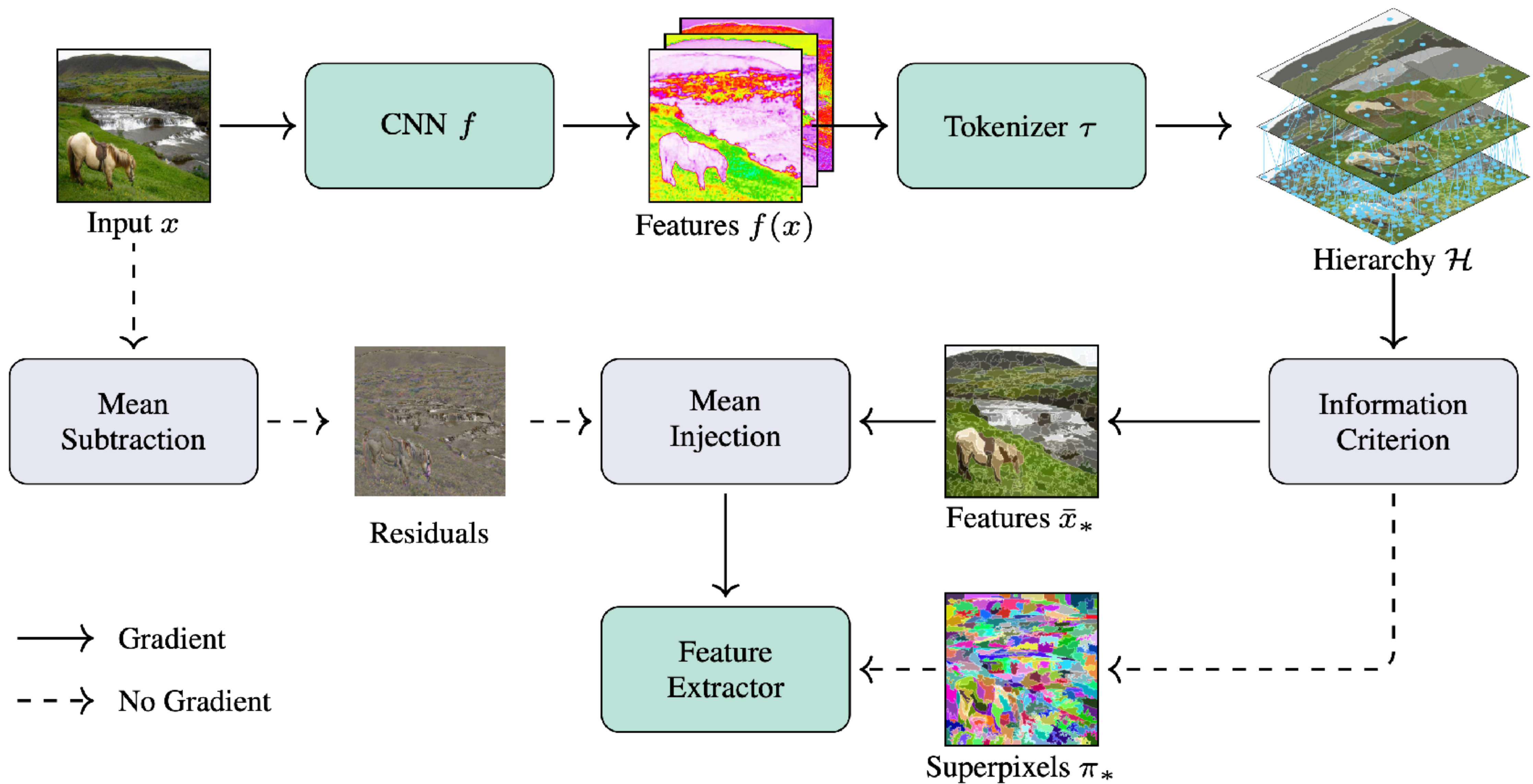
An army of aardvarks awkwardly asked their aunties for advice on anteaters.

An army of aardvarks awkwardly asked their aunties for advice on anteaters.

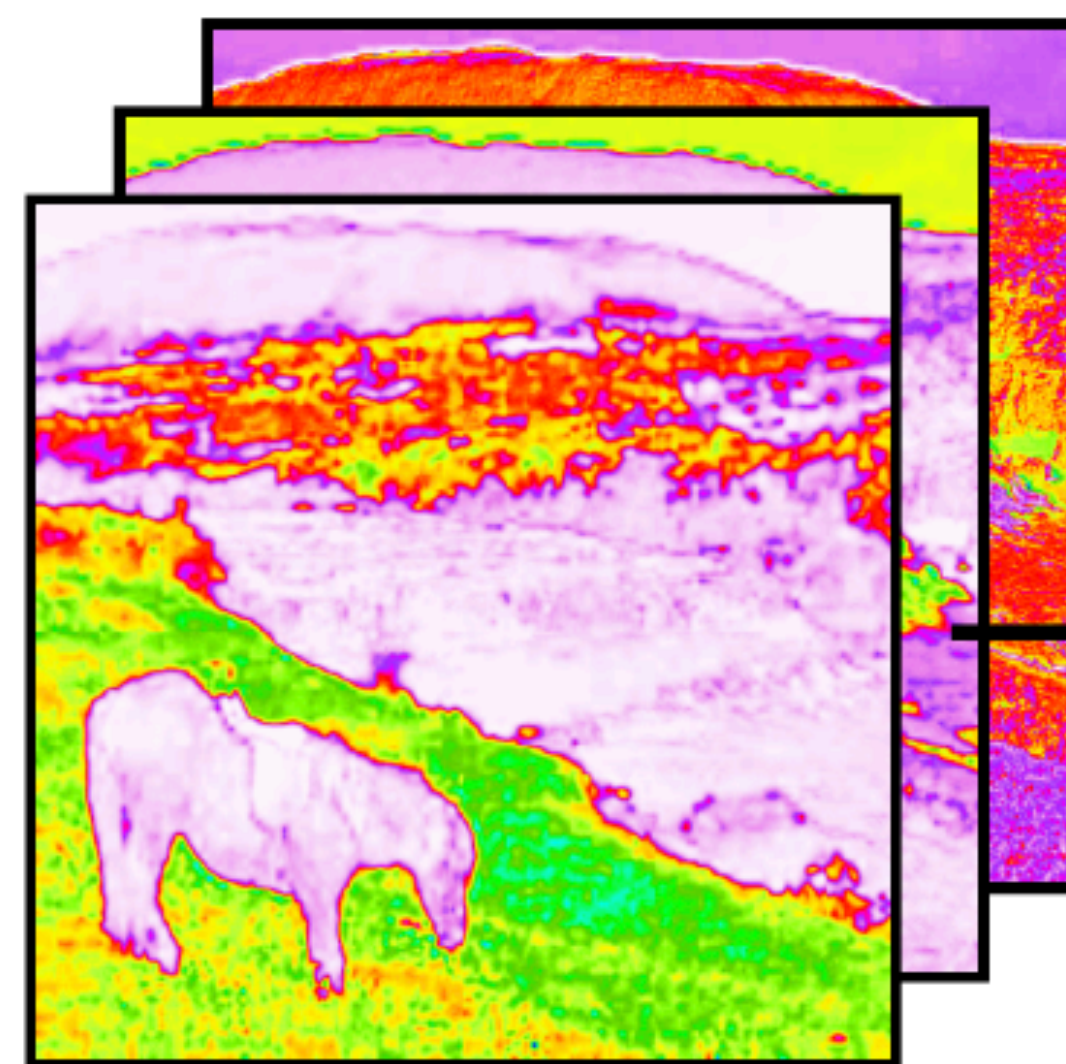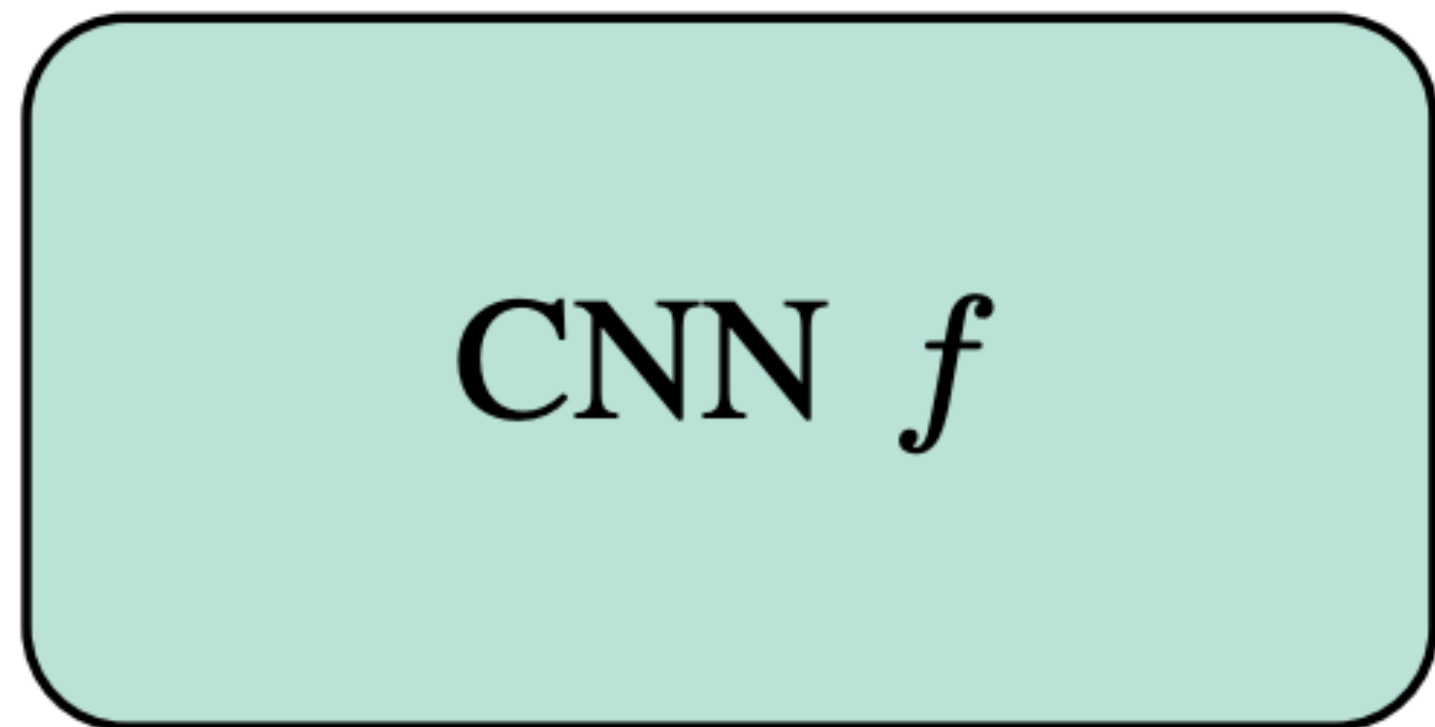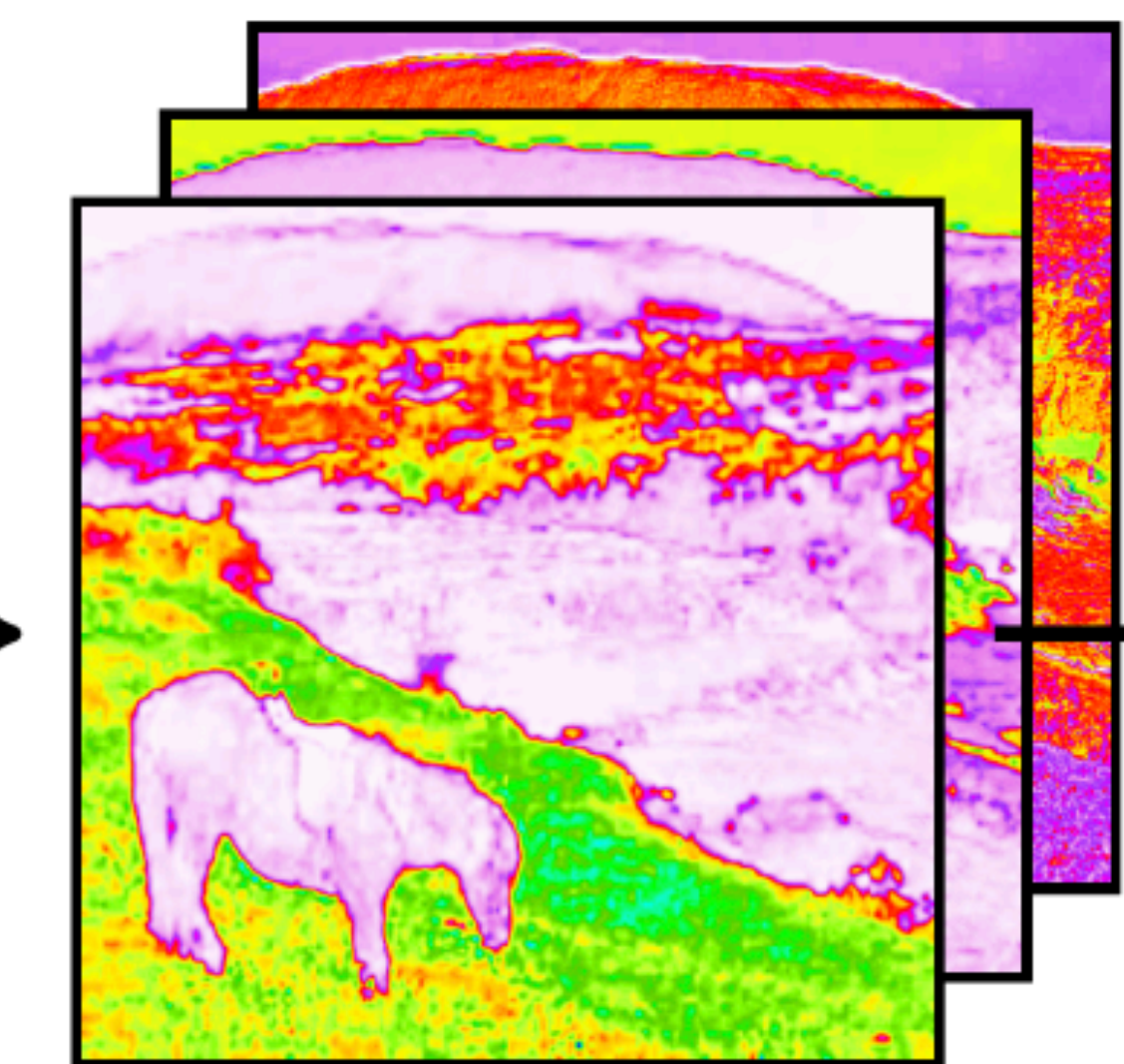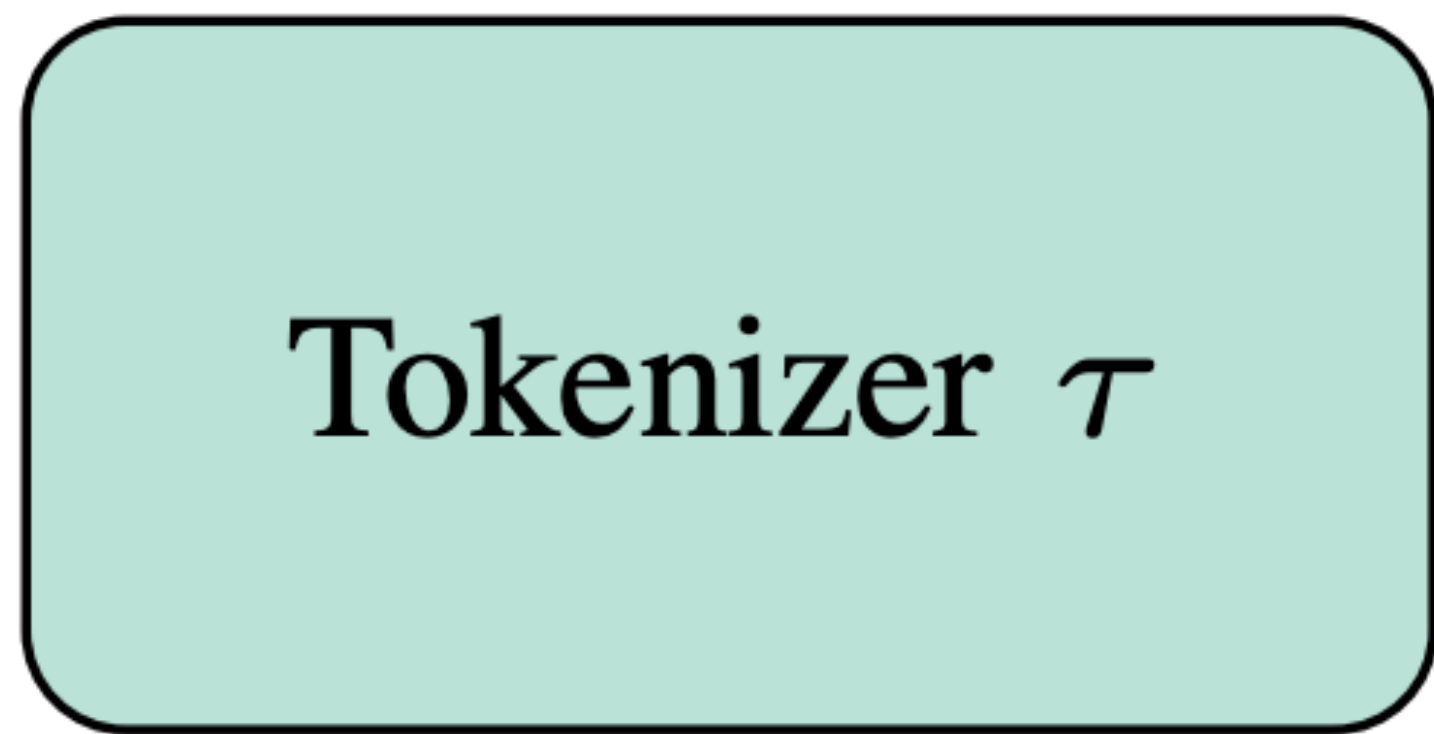An army of aardvarks awkwardly asked their aunties for advice on anteaters.

Input $x$

CNN $f$

Features $f(x)$

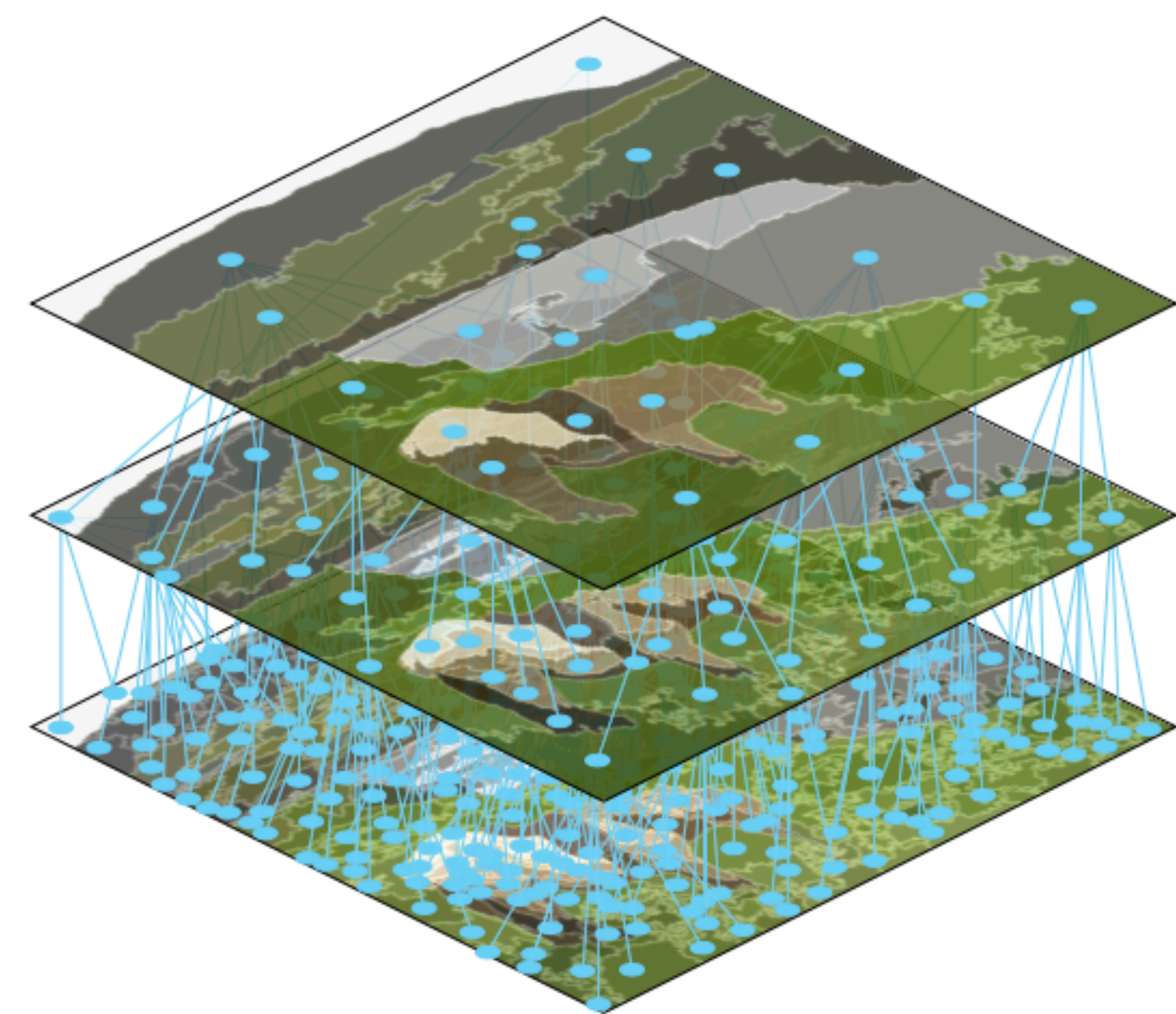Features $f(x)$

Tokenizer $\tau$

Hierarchy $\mathcal{H}$

Mean Injection

Features $\bar{x}_*$

Information Criterion

Feature Extractor

Superpixels $\pi_*$

Residuals

Mean Injection

Feature Extractor

Features $\bar{x}_*$

Superpixels $\pi_*$

**Feature Extractor**

$$F(S) = (\;\blacksquare_{M^+} + \lambda\;\blacksquare_{M^-}\;) \odot \hat{x}(S)$$
$$+ (1 - \lambda)\;\blacksquare_{M^-} \odot \beta$$

**Figure 1:** Comparing spatial granularity in visual tokenizers. $\partial$HT (right) provides an end-to-end learnable framework for multi-scale tokenization. We provide more examples in Figure E.6.
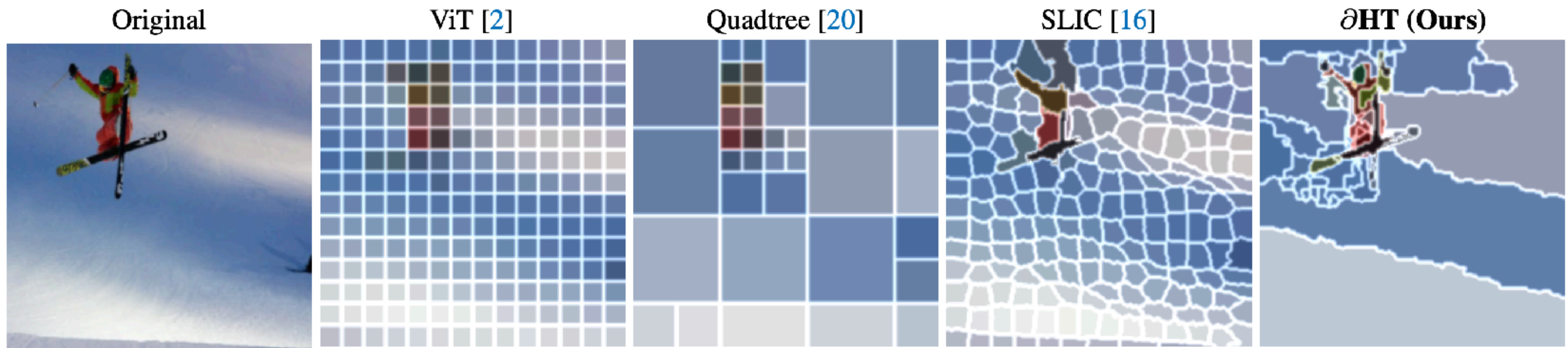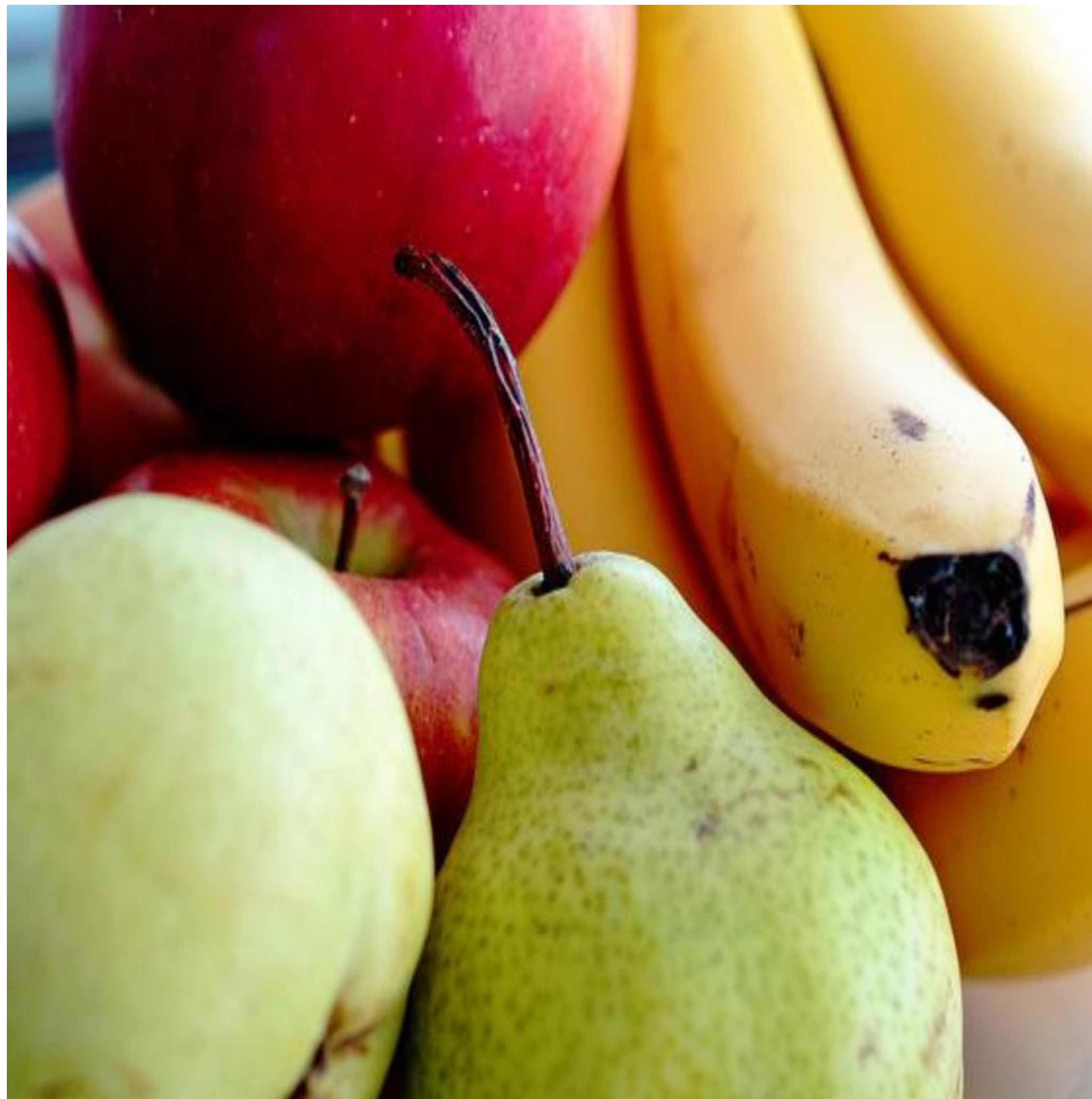
**Figure E.1:** Segmentation examples for $\partial$HT over ADE20k, showing fine grained segmentation labels only using a simple MLP head without upscaling. *Top*: Original images (512 × 512). *Middle*: Annotated target images. *Bottom*: Predicted labels from $\partial$HT.

Original



Vectorized

Original



Vectorized

Original

Vectorized

Original

Vectorized

# Thank you for your attention!

**Marius Aasan - University of Oslo | SFI Visual Intelligence**