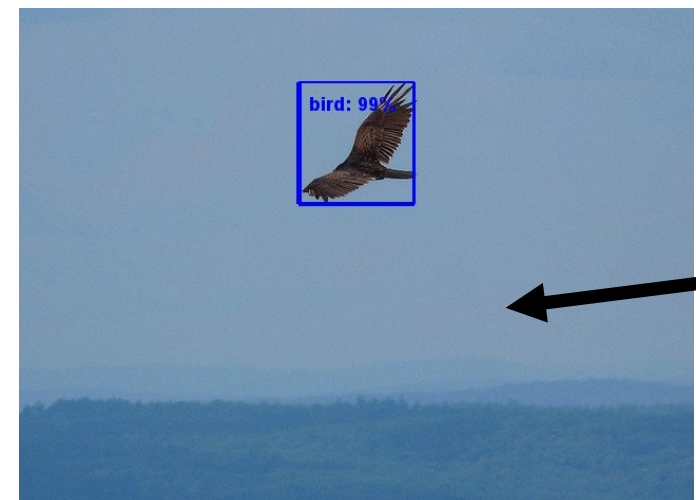


# Spatially Adaptive Computation Time for Residual Networks

Michael Figurnov<sup>1</sup> Maxwell D. Collins<sup>2</sup> Yukun Zhu<sup>2</sup> Li Zhang<sup>2</sup> Jonathan Huang<sup>2</sup> Dmitry Vetrov<sup>1,3</sup> Ruslan Salakhutdinov<sup>4</sup>

<sup>1</sup>National Research University Higher School of Economics <sup>2</sup>Google Inc. <sup>3</sup>Yandex <sup>4</sup>Carnegie Mellon University

## Motivation



Do we really need to put these blue pixels through 100 layers to detect the bird?

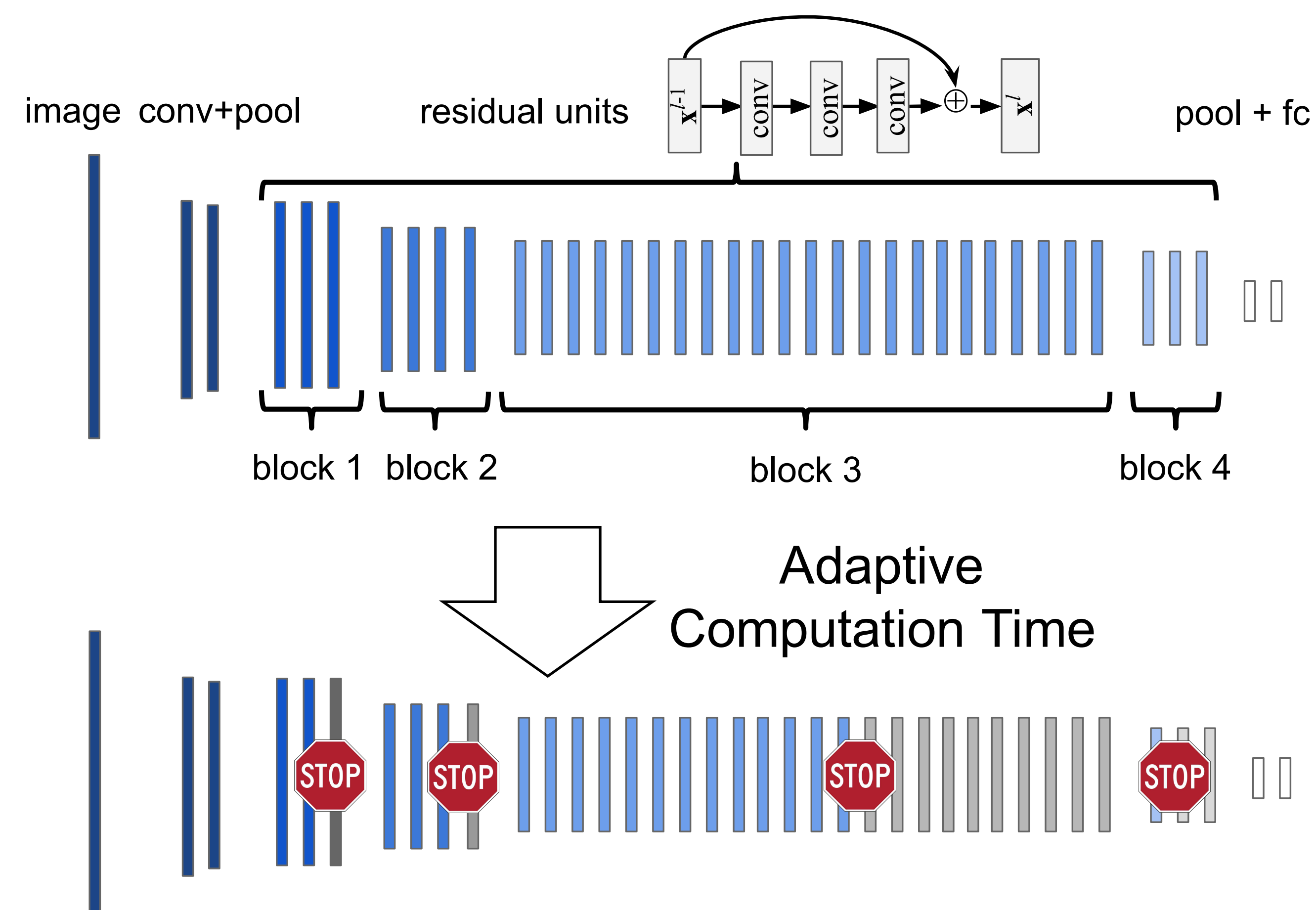
Can we make networks **faster** and **more interpretable** by adapting the amount of computation?

## Contribution

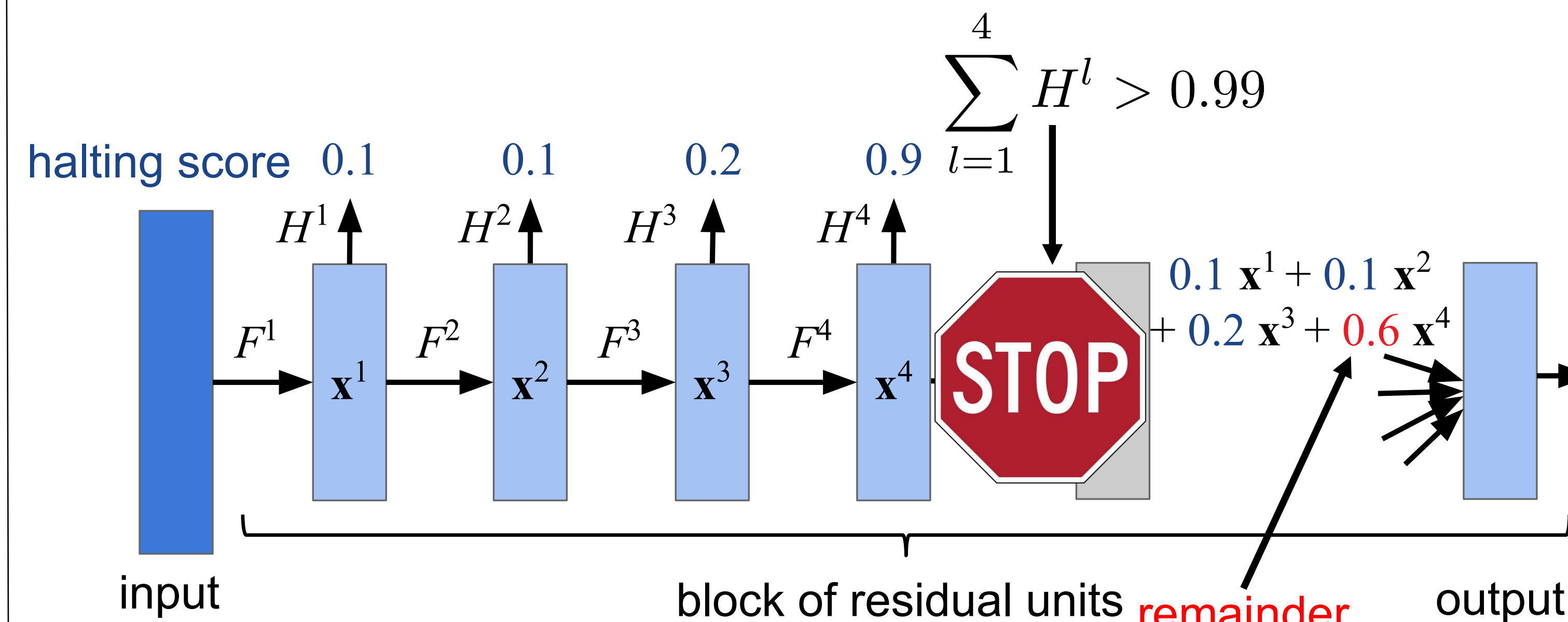
A novel mechanism that

- adapts the computation **spatially**
- is end-to-end trainable
- provides introspection
- scales to ImageNet and COCO
- is problem-agnostic

## Residual Network (ResNet)



## Adaptive Computation Time (ACT)

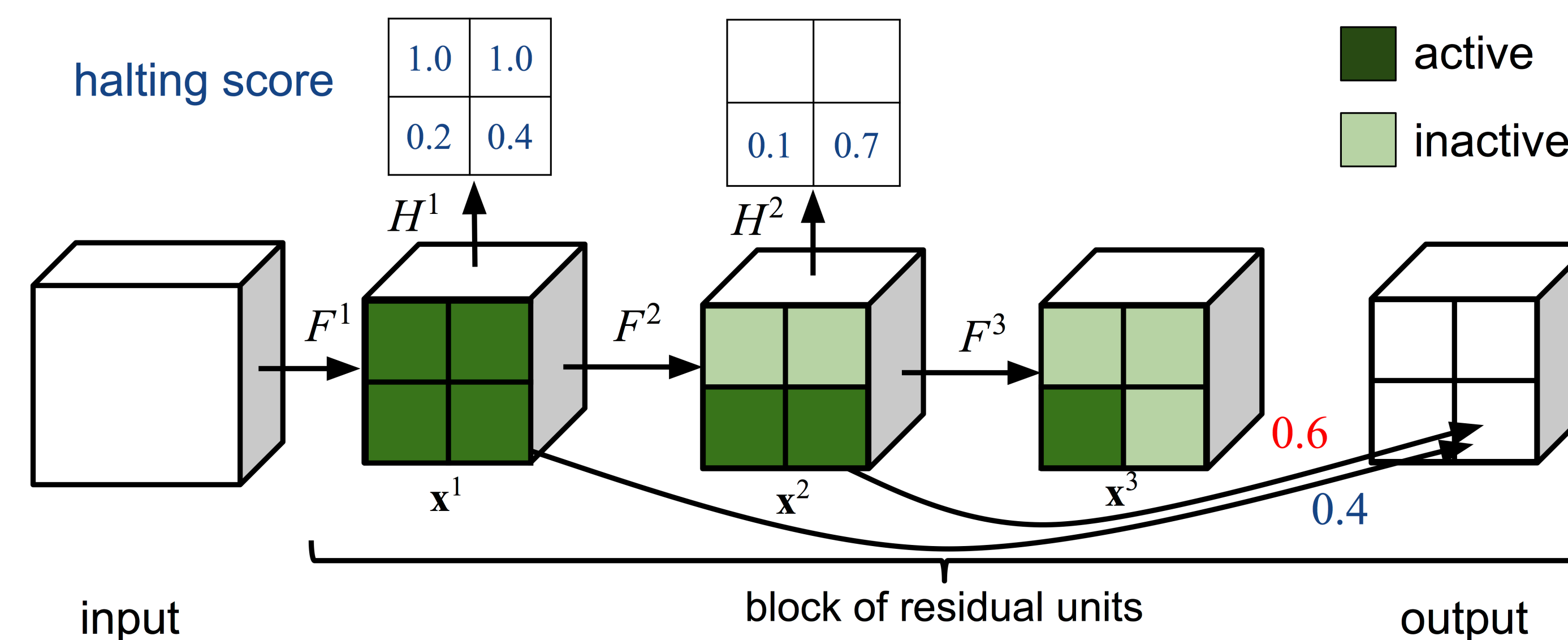


$$\text{Loss function } \mathcal{L}' = \mathcal{L} + \tau(N + R), \quad R = 1 - \sum_{l=1}^{N-1} H^l$$

original loss      number of residual units      ponder cost (close to  $N$ , but differentiable)

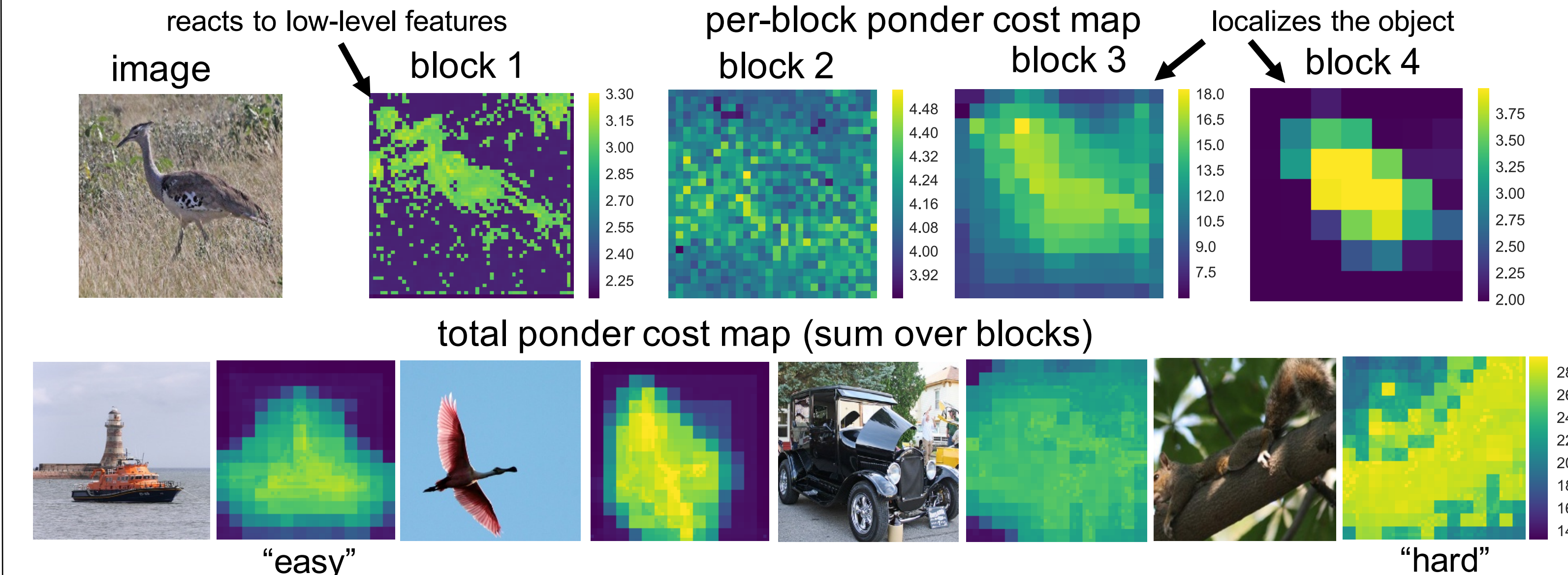
## Spatially Adaptive Computation Time (SACT)

Apply ACT “convolutionally” to every spatial position of residual network’s block

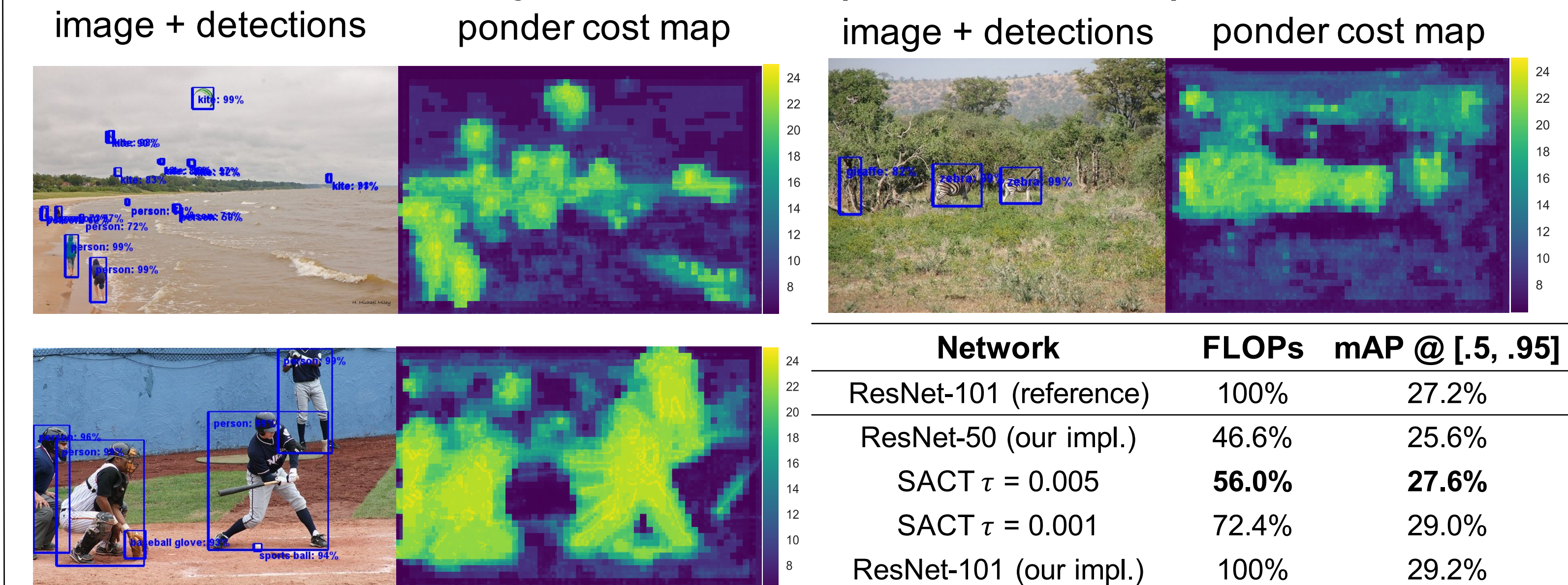


Residual units are evaluated only in the active positions!

## Image classification (ImageNet dataset)

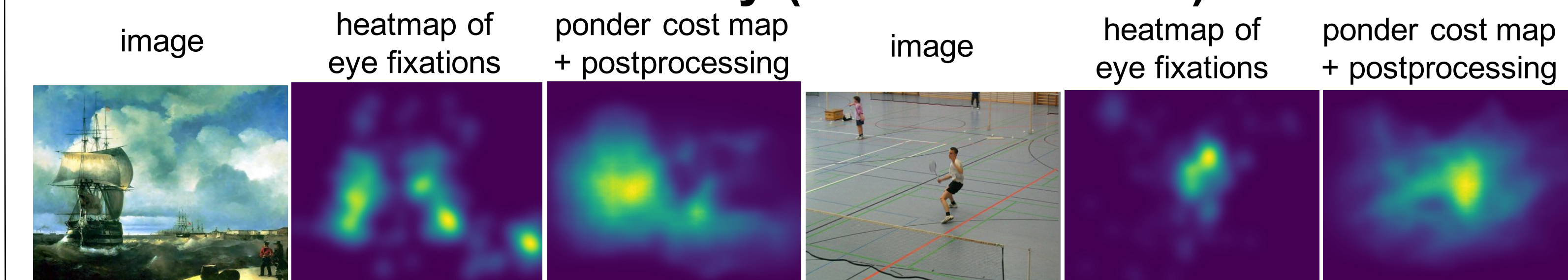


## Object detection (COCO dataset)



SACT has better speed-mAP trade-off compared to reducing the depth of ResNet

## Visual saliency (cat2000 dataset)



Ponder cost maps tend to highlight salient objects, even though the SACT model is not trained on cat2000!

Code:  
github.com/mfigurnov/sact