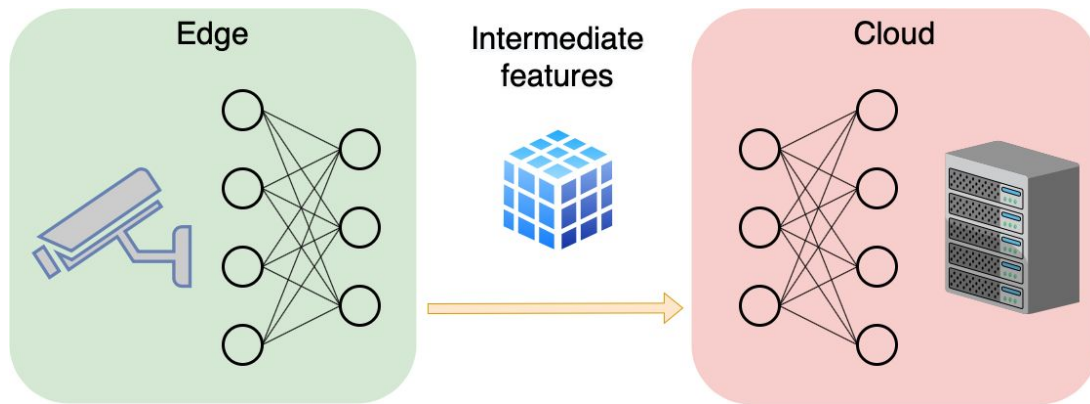


Latent Space Motion Analysis for Collaborative Intelligence

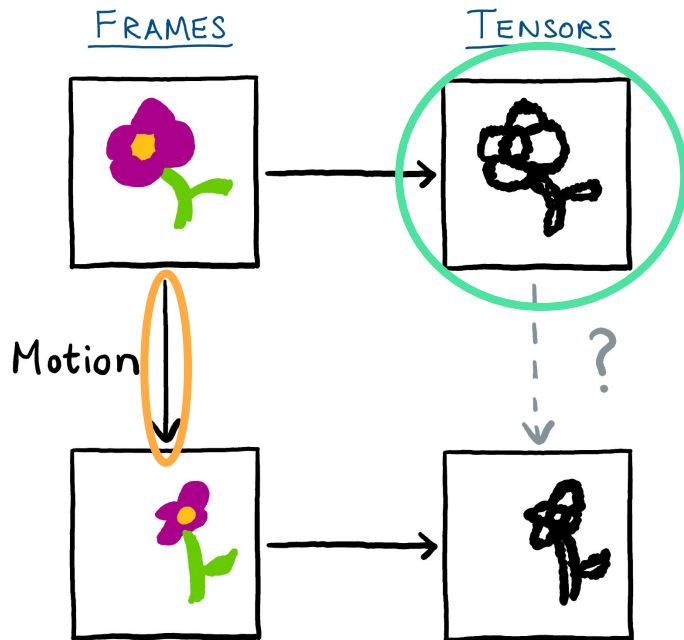
Mateen Ulhaq
Ivan V. Bajić



Problem statement

Processing an input frame through part of a deep model yields a feature tensor.

Question: Given a **reference tensor** and the **motion** between successive input frames, can we determine the motion between successive tensors?





Observation: tensor motion \approx rescaled input motion

Three types of operations in CNNs:

- Convolution
- Pointwise activation (e.g. ReLU)
- Pooling (e.g. max)

Show: motion after operation \approx motion before operation

\Rightarrow motion after multiple operations \approx motion in input

Optical flow

$$\text{2D} \quad \frac{\partial I}{\partial x} v_x + \frac{\partial I}{\partial y} v_y + \frac{\partial I}{\partial t} = 0$$

$I(x, y)$
pixel intensity

$$\text{1D} \quad \frac{\partial I}{\partial x} v + \frac{\partial I}{\partial t} = 0$$

(v_x, v_y)
motion vectors

Convolution

Input
optical flow:

$$\frac{\partial I}{\partial x} v + \frac{\partial I}{\partial t} = 0$$

After
convolution:

$$\frac{\partial}{\partial x} [f * I] \tilde{v} + \frac{\partial}{\partial t} [f * I] = 0$$

f = conv filter
 \tilde{v} = post-conv motion

By commutativity of
linear operators:

$$f * \left(\frac{\partial I}{\partial x} \tilde{v} + \frac{\partial I}{\partial t} \right) = 0$$

one solution to this equation is

$$\tilde{v} = v$$

Pointwise activation (e.g. ReLU)

Input

optical flow:

$$\frac{\partial I}{\partial x} v + \frac{\partial I}{\partial t} = 0$$

After activation:

$$\frac{\partial}{\partial x} [\sigma(I)] \tilde{v} + \frac{\partial}{\partial t} [\sigma(I)] = 0$$

σ = activation
 \tilde{v} = post-activation
motion

By chain rule:

$$\sigma'(I) \cdot \left(\frac{\partial I}{\partial x} \tilde{v} + \frac{\partial I}{\partial t} \right) = 0$$

one solution to this equation is

$$\tilde{v} = v$$

Max pool

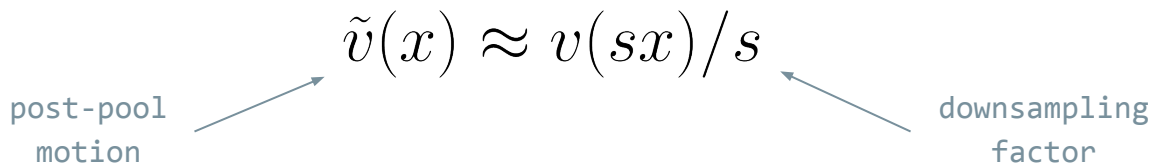
Input
optical flow:

$$\frac{\partial I}{\partial x}v + \frac{\partial I}{\partial t} = 0$$

Derivation:

[described in paper]

Result:



The diagram shows the equation $\tilde{v}(x) \approx v(sx)/s$ with two arrows pointing towards it. One arrow originates from the text 'post-pool motion' and points to the left side of the equation. The other arrow originates from the text 'downsampling factor' and points to the right side of the equation.

$$\tilde{v}(x) \approx v(sx)/s$$

post-pool
motion

downsampling
factor

Tensor reconstruction experiments

- Reconstruct current tensor by applying the rescaled motion calculated between previous and current frame.

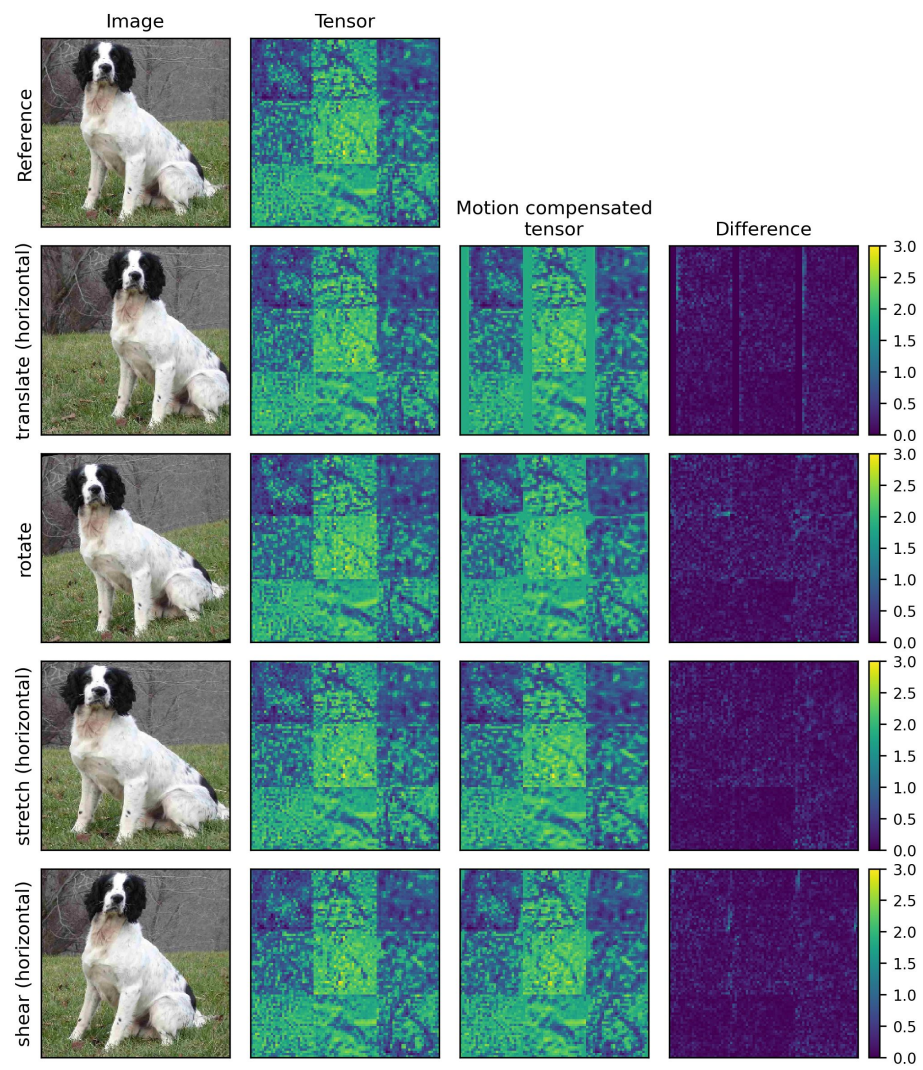
pool layers

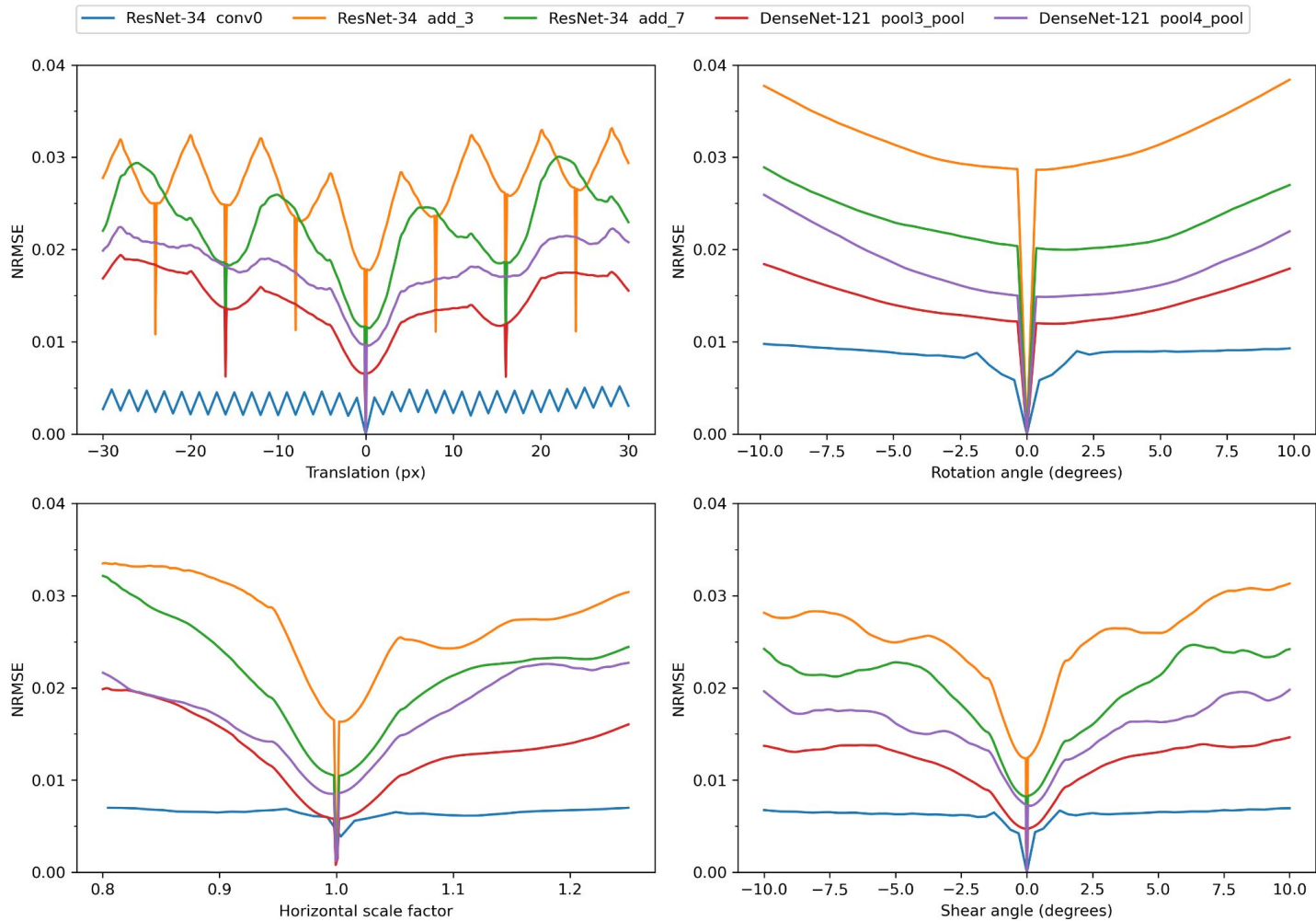


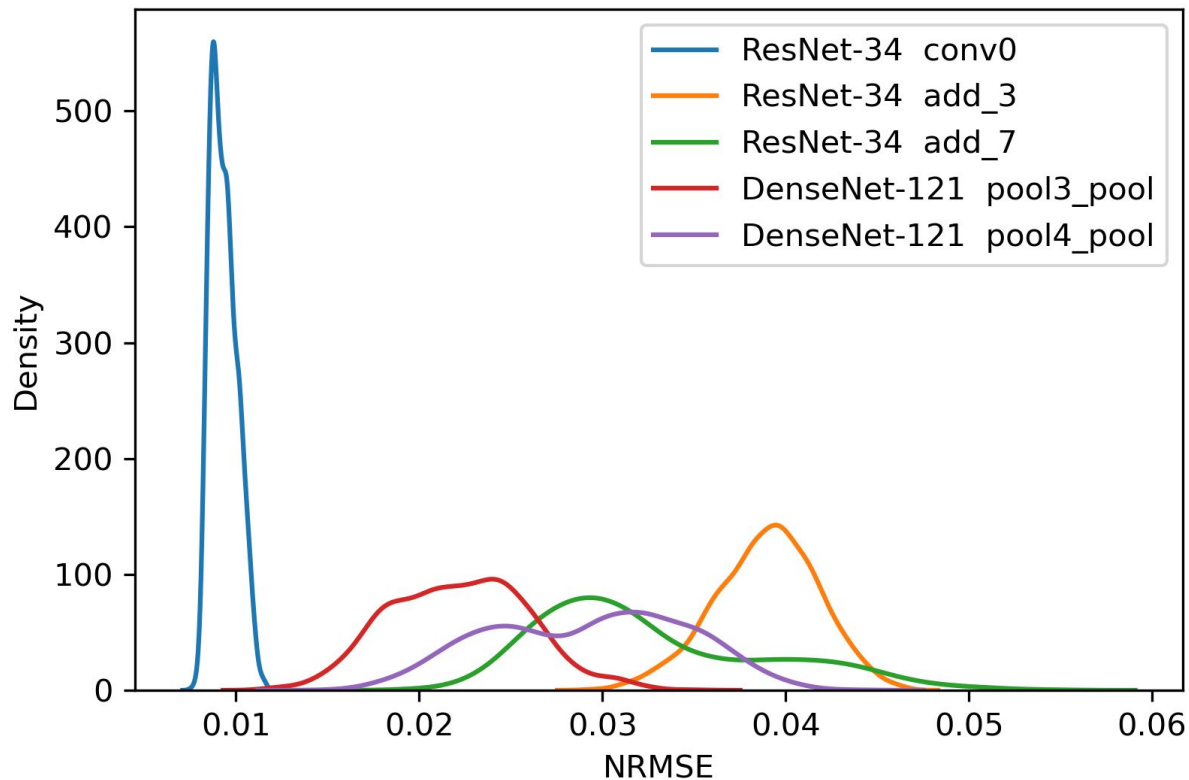
$$\tilde{v}(x, y) \approx v(s^k x, s^k y) / s^k$$

- Calculate normalized root mean square error (NRMSE) between predicted tensor and ground truth tensor.

$$\text{NRMSE} = \frac{1}{R} \sqrt{\frac{1}{N} \sum_{i=1}^N (p_i - a_i)^2}$$







Affine transformation
composed of:

- x, y translation (± 32 px)
- x, y scaling ($0.95 - 1.05x$)
- x, y shearing ($\pm 5^\circ$)
- rotation ($\pm 10^\circ$)

NRMSE of 0.04 roughly
corresponds to 28 dB PSNR in
traditional video motion.

$$\frac{\partial}{\partial x}[f * I]\tilde{v} + \frac{\partial}{\partial t}[f * I] = 0$$

$$f * \left(\frac{\partial I}{\partial x} \tilde{v} + \frac{\partial I}{\partial t} \right) = 0$$



one solution to this equation is

$$\tilde{v} = v$$

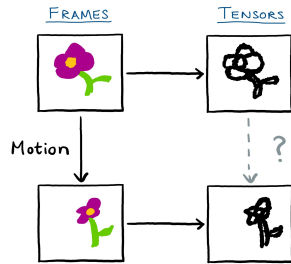
Conclusion

- Derived a simple approximate relationship of motion within the tensor channels by analyzing typical operations in CNNs.

$$\tilde{v}(x, y) \approx v(s^k x, s^k y) / s^k$$

- The reconstruction error for small affine transformations within the input is 4% (NRMSE).

Thank you



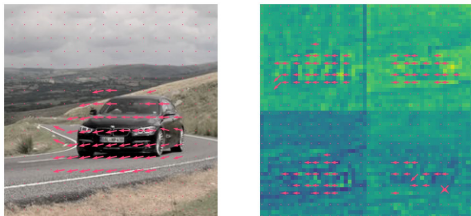
Question: What is the relationship between input-space motion and latent-space motion?

By analyzing optical flow before and after typical CNN operations, we show:

$$\tilde{v}(x, y) \approx v(s^k x, s^k y) / s^k$$

input motion

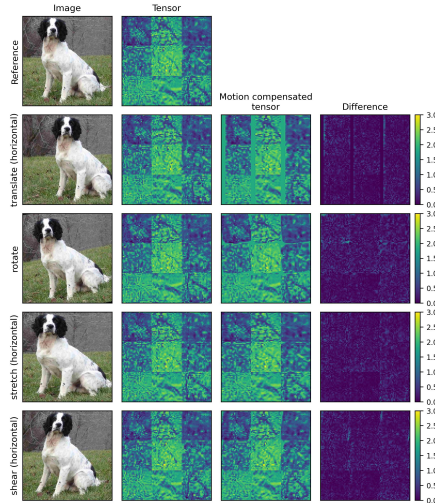
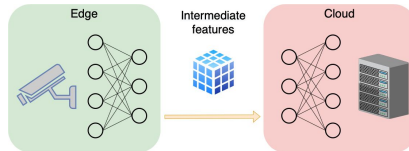
stride # pooling layers



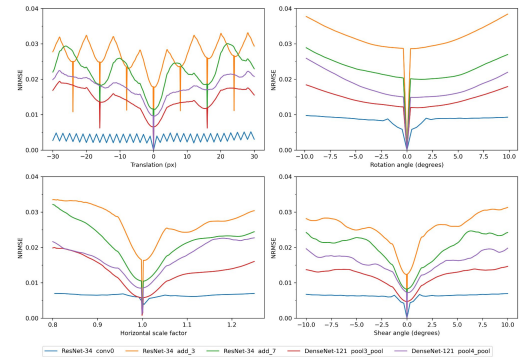
Estimates of input motion (left) and several channels from the output of ResNet-34's add_3 layer (right).

Latent Space Motion Analysis for Collaborative Intelligence

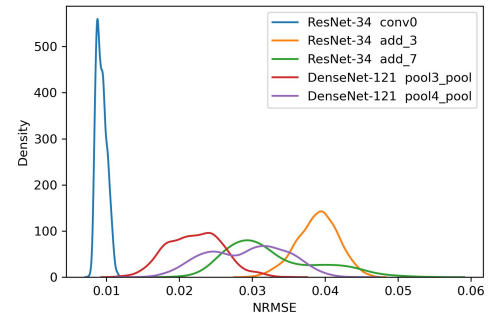
Mateen Ulhaq and Ivan V. Bajić



Various transformations applied to reference image. The output tensors of ResNet-34's add_3 layer are reliably predicted from only the reference tensor and known input-space transformation.



NRMSE for translation (top-left), rotation (top-right), scaling (bottom-left), and shear (bottom-right). For translation, NRMSE local minima occur when the input-space shifts correspond to integer latent-space shifts.



NRMSE histogram for reconstruction of affine-transformed inputs with translation (± 32 px), scaling ($0.95x - 1.05x$), shearing ($\pm 5^\circ$), rotation ($\pm 10^\circ$). NRMSE of $0.04 \approx 28$ dB PSNR.

Max pool

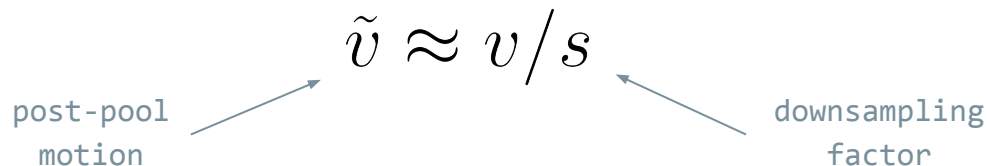
Input
optical flow:

$$\frac{\partial I}{\partial x} v + \frac{\partial I}{\partial t} = 0$$

Derivation:

[described in paper]

Result:



The diagram shows the equation $\tilde{v} \approx v/s$ in the center. Two blue arrows point towards this equation. The arrow from the left is labeled "post-pool motion" and the arrow from the right is labeled "downsampling factor".

$$\tilde{v} \approx v/s$$

post-pool motion

downsampling factor

Tensor reconstruction experiments

- Reconstruct current tensor by applying the rescaled motion calculated between previous and current frame.

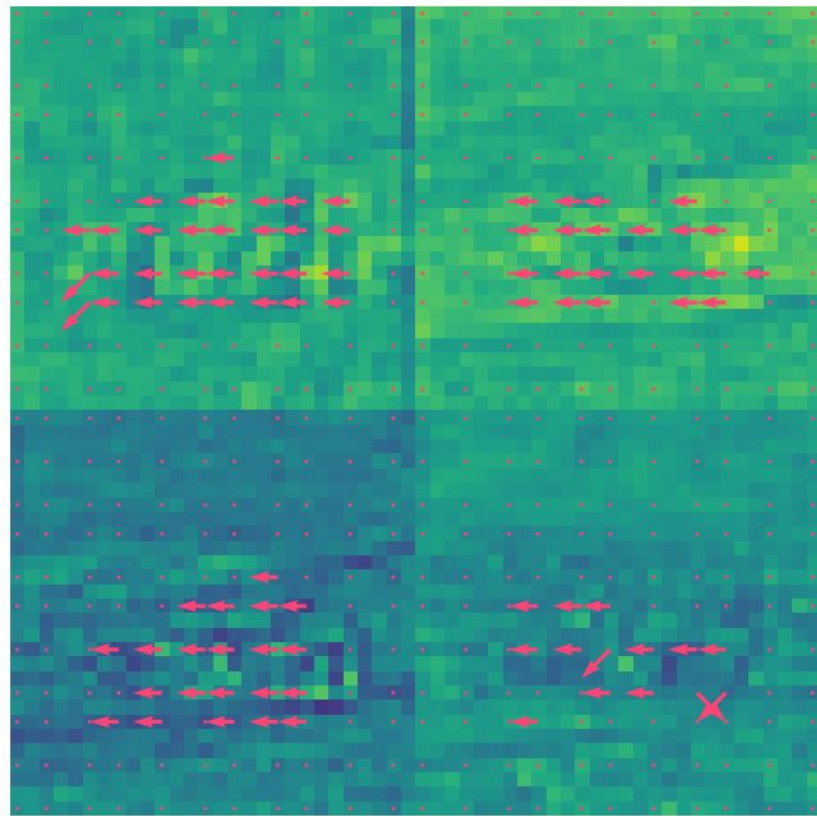
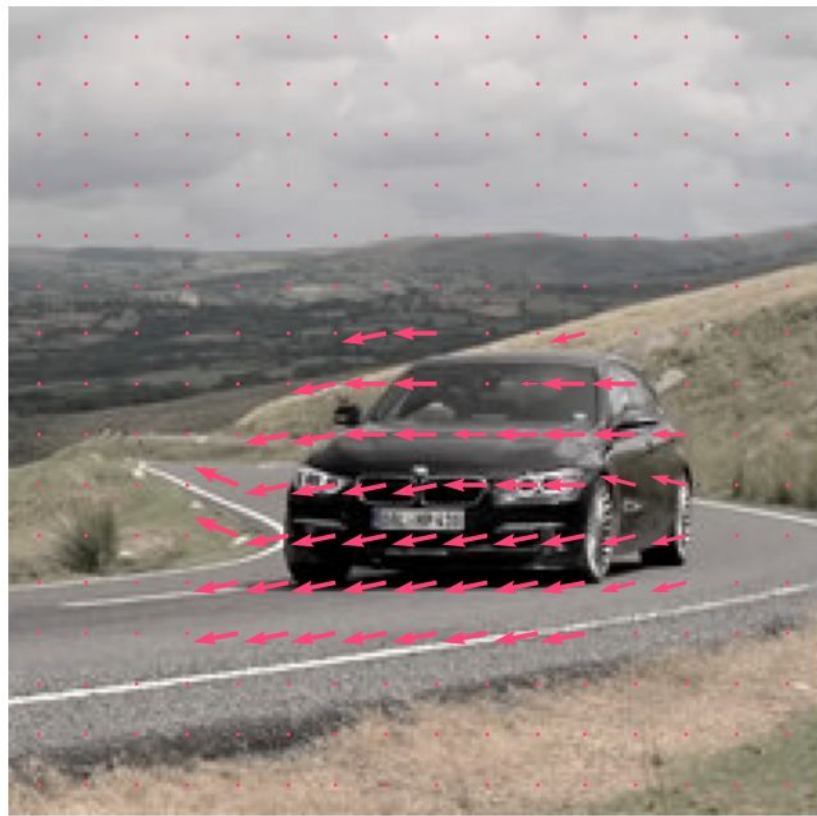
pool layers

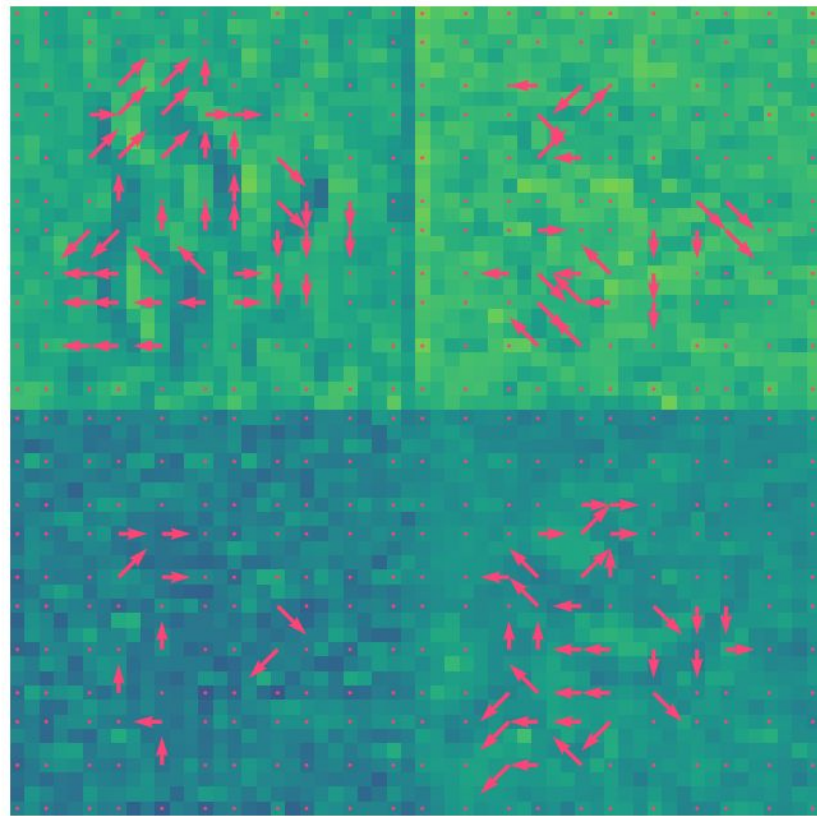


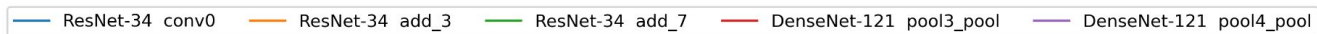
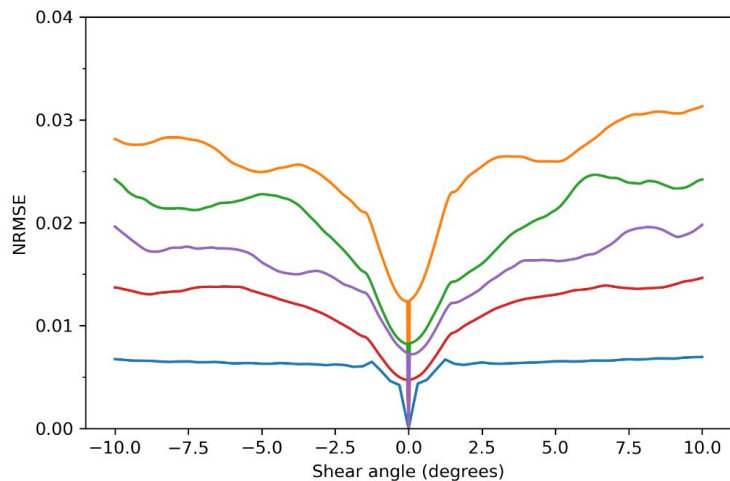
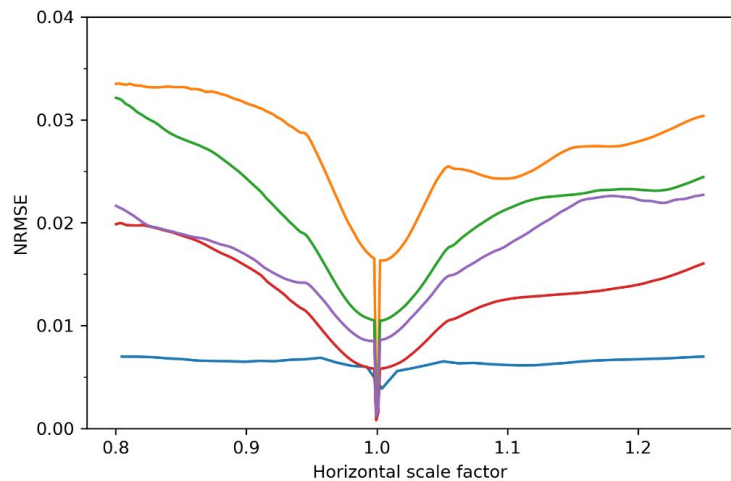
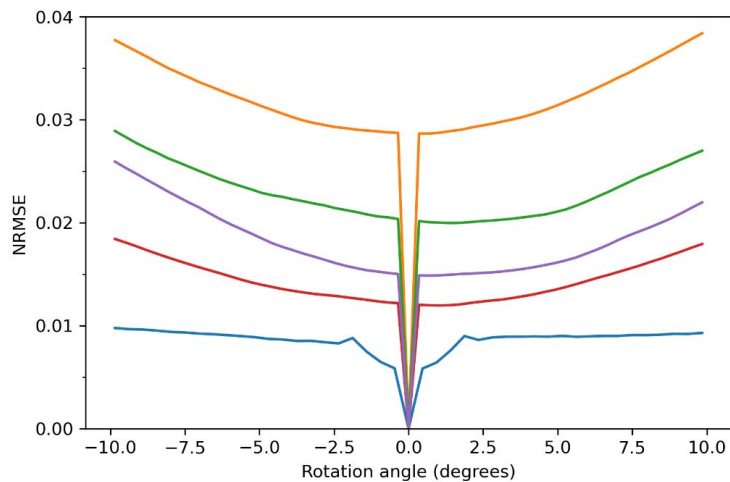
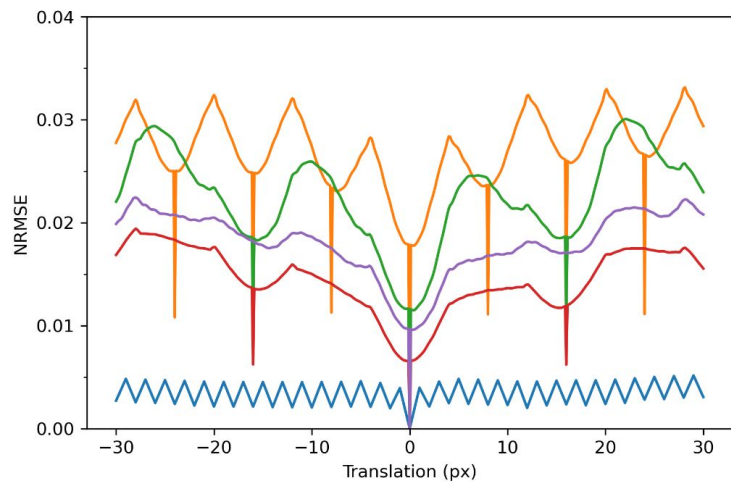
$$\tilde{v} \approx v / s^k$$

- Calculate normalized root mean square error (NRMSE) between predicted tensor and ground truth tensor.

$$\text{NRMSE} = \frac{1}{R} \sqrt{\frac{1}{N} \sum_{i=1}^N (p_i - a_i)^2}$$







Shared inference

Key idea: less data sent over network

Versus cloud-only inference:

- Save bandwidth
- Save device energy
- Reduce inference times

Versus edge-only inference:

- Bigger models
- Reduce resource usage
- Reduce inference times

