



# Machine Learning in Graphics and Vision

**- Global Illumination & ML -**

SoSe 2018

Hendrik Lensch





# Previous Work

	Non-linear Playback	Appearance Changes	Temporal Alignment	Handling Detours
a Evangelidis et al. (2011)				
b Torii et al. (2015)				
c Douze et al. (2016)				
d				



## Efficiency: Sparse Feature-Extraction (260h Data)

Approach	FPS	Total Time
Evangelidis et al. (2011)	0.11	0.81 years (~42 weeks)
Evangelidis et al. (2013)	0.11	0.81 years (~42 weeks)
Wang et al. (2014)	3.77	8.6 days

Ours	140	5.57 hours
------	-----	------------

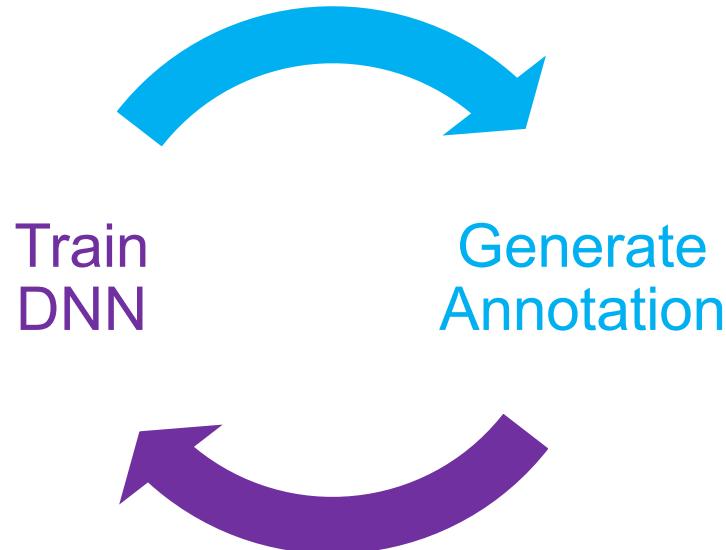
x1272 speed up

Ours (12 GPUS)	12 * 140	28 minutes
----------------	----------	------------



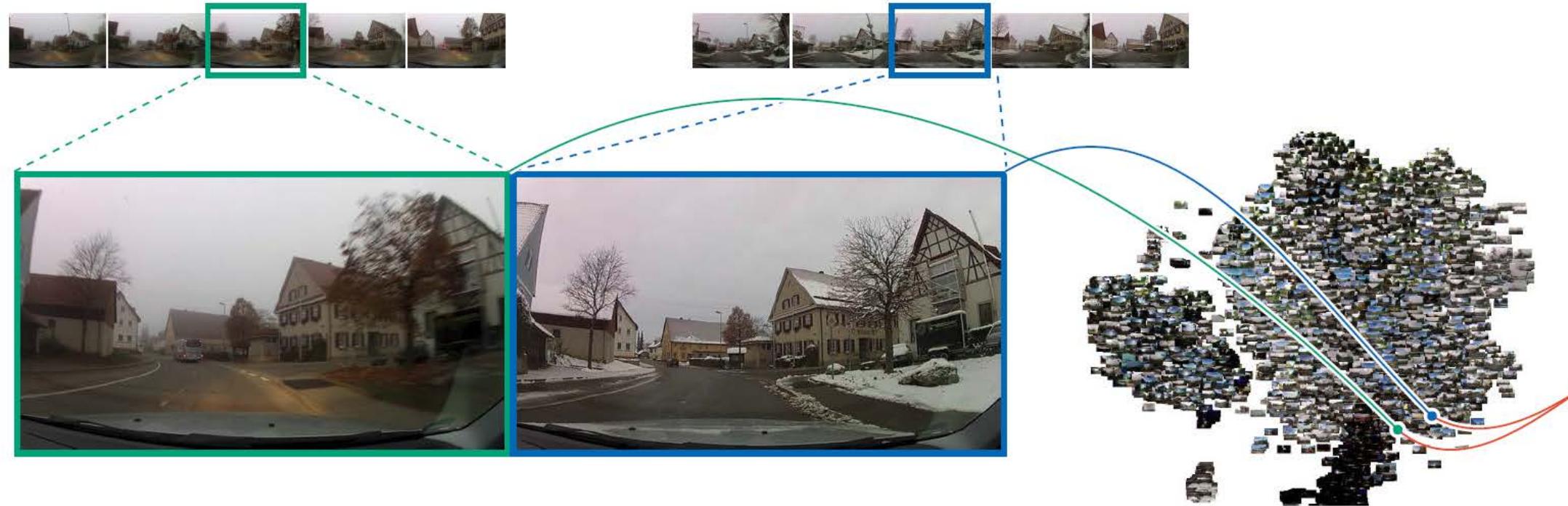
# Data-Driven Approaches without Annotations

- Chicken-Egg-Problem





# Encoding each Frame

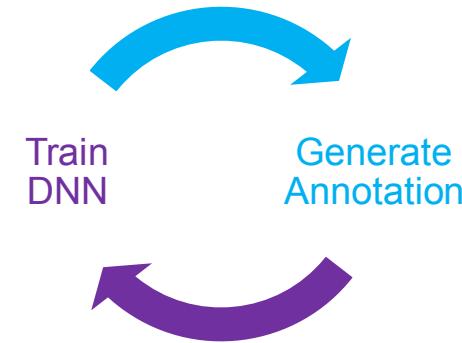
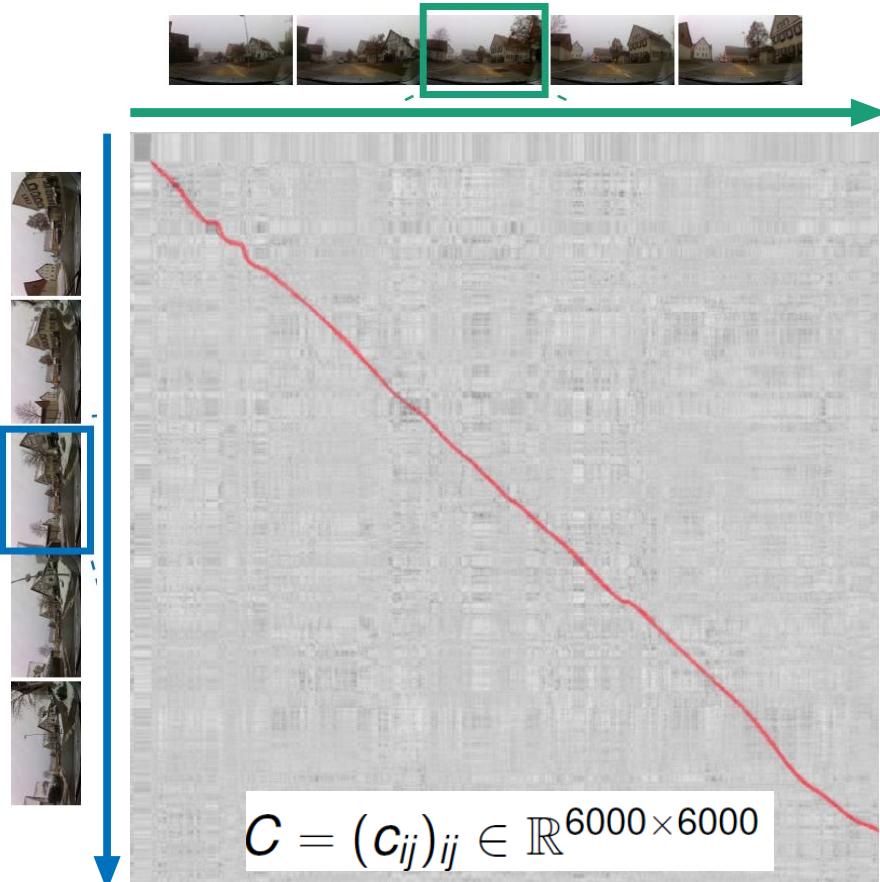


- Visualization (t-SNE): mapping 1024 dimensions to 2 dimensions
- Euclidean distance between two frames in high-dimensional encoding-space represents their similarity.



## Phase: Generate Annotation

- Assume trained Deep Neural Network is given



- pair-wise distances  $C'$
- low-rank approximation
- de-correlated distances  $C$
- find synchronization path



# Finding a Synchronization Path

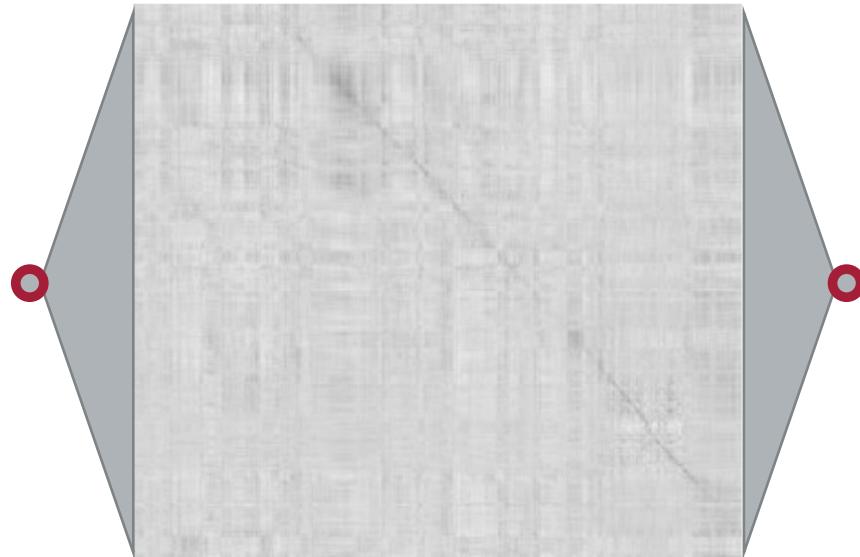
- Dijkstra's algorithm finds shortest path between two given nodes.
- We do not know the start or end node!





# Finding a Synchronization Path

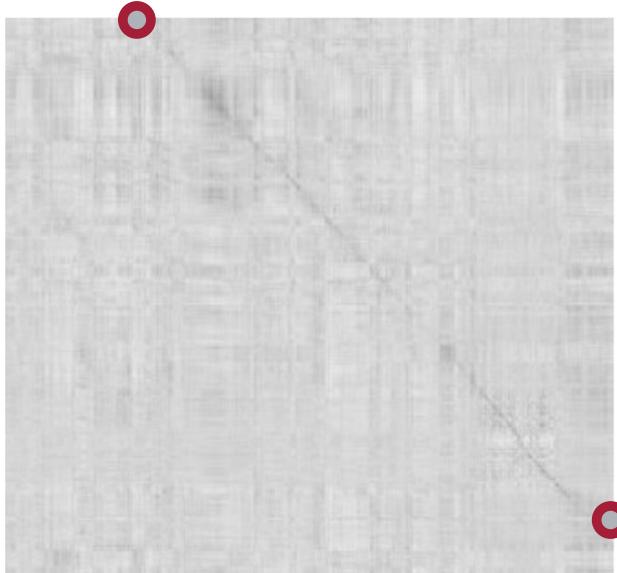
- Dijkstra's algorithm finds shortest path between two given nodes.
- We do not know the start or end node!





# Finding a Synchronization Path

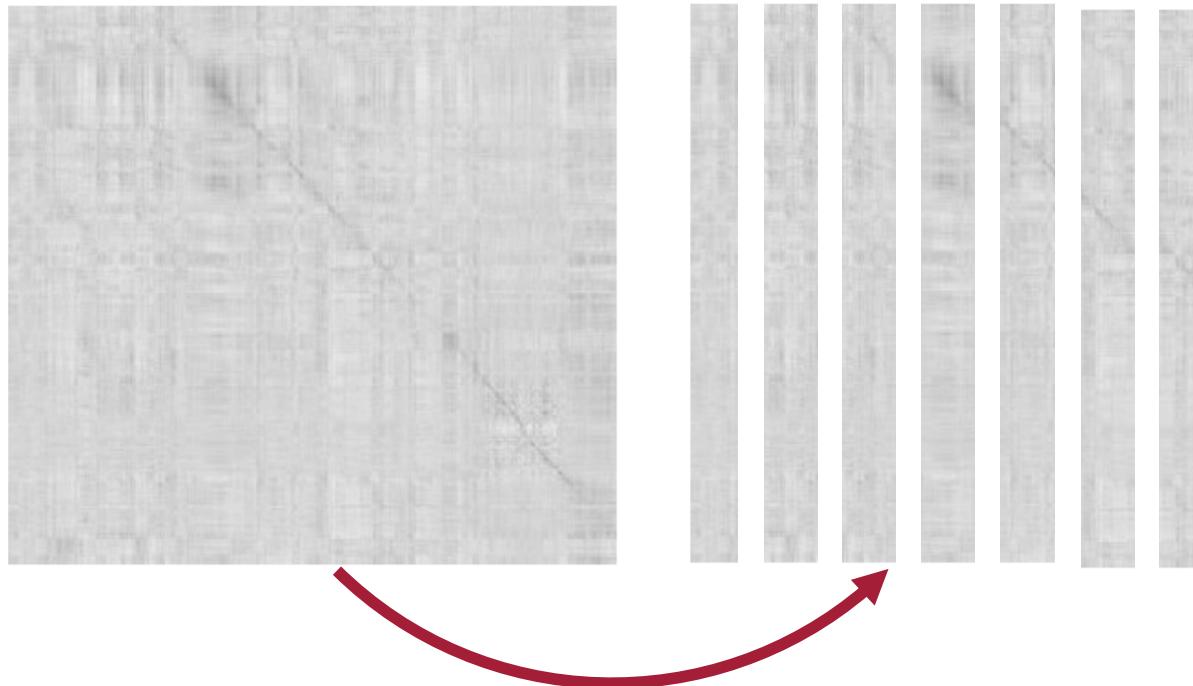
- Dijkstra's algorithm finds shortest path between two given nodes.
- We do not know the start or end node!
- We even do not know if there is any path!





# Finding a Synchronization Path

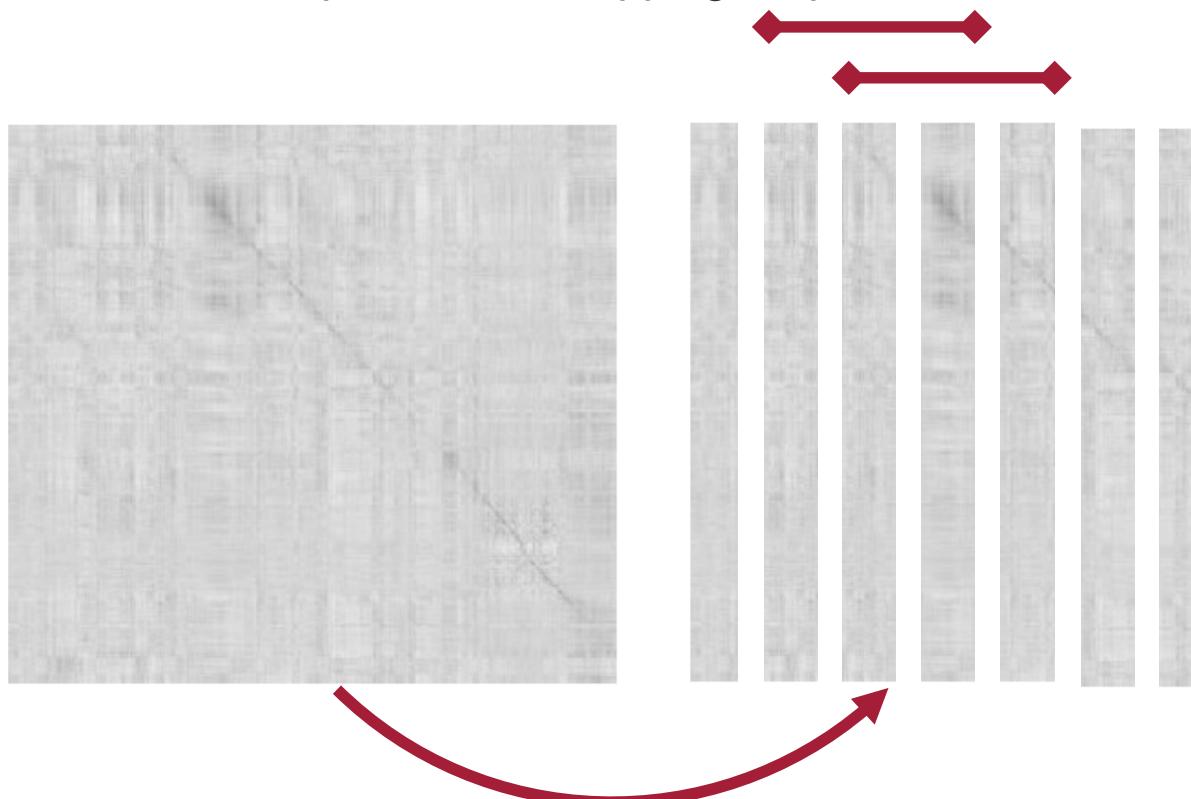
- Dijkstra's algorithm finds shortest path
- Partition into strips





# Finding a Synchronization Path

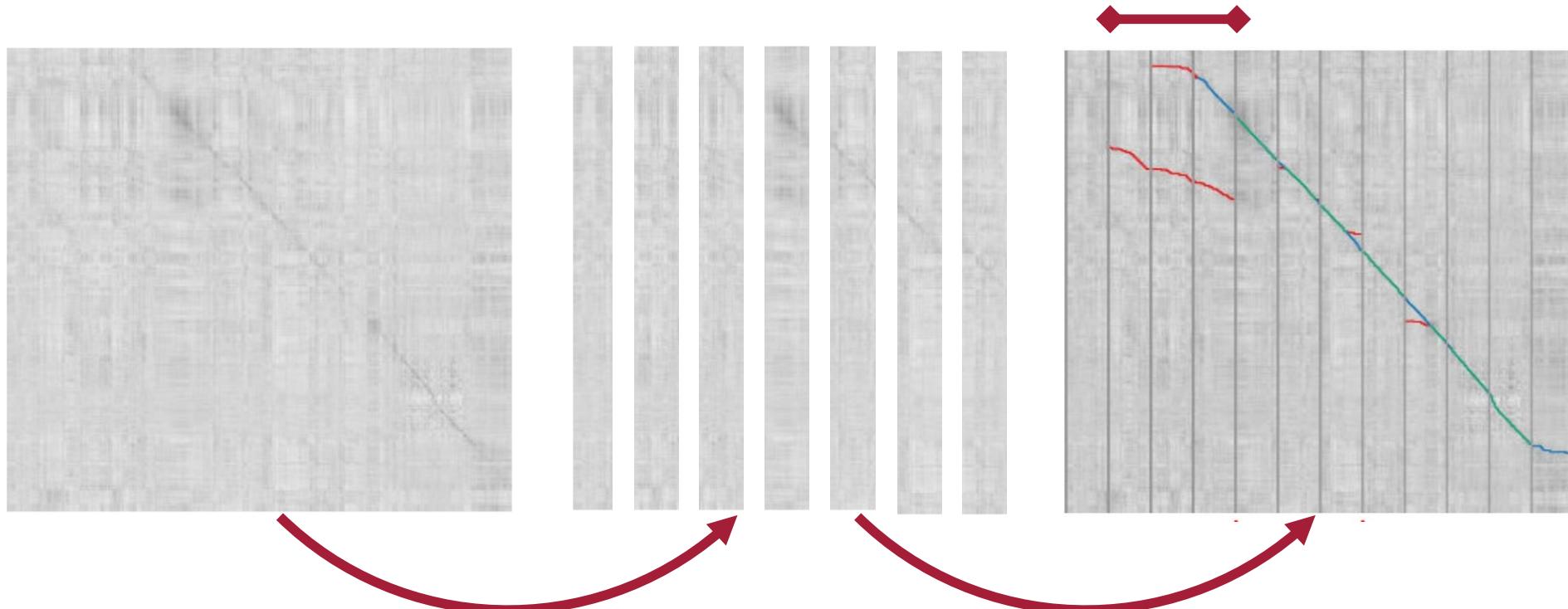
- Dijkstra's algorithm finds shortest path
- Partition into strips
- Find shortest path for overlapping strips





# Finding a Synchronization Path

- Dijkstra's algorithm finds shortest path
- Partition into strips
- Find shortest path for overlapping strips
- Judge consistency in overlap regions

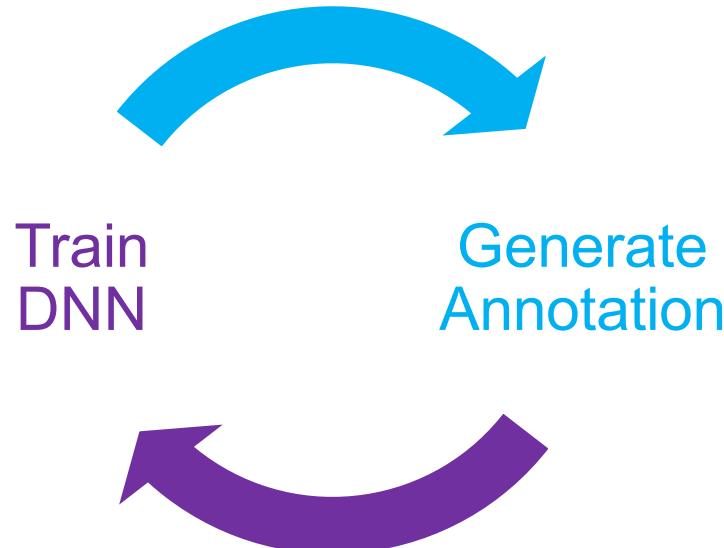






# Data-Driven Approaches without Annotations

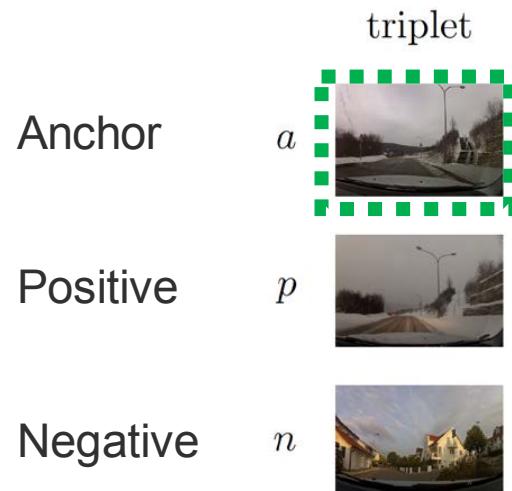
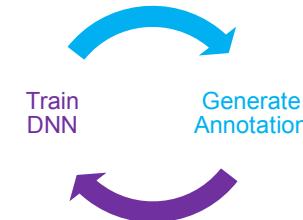
- Chicken-Egg-Problem





## Phase: Train DNN

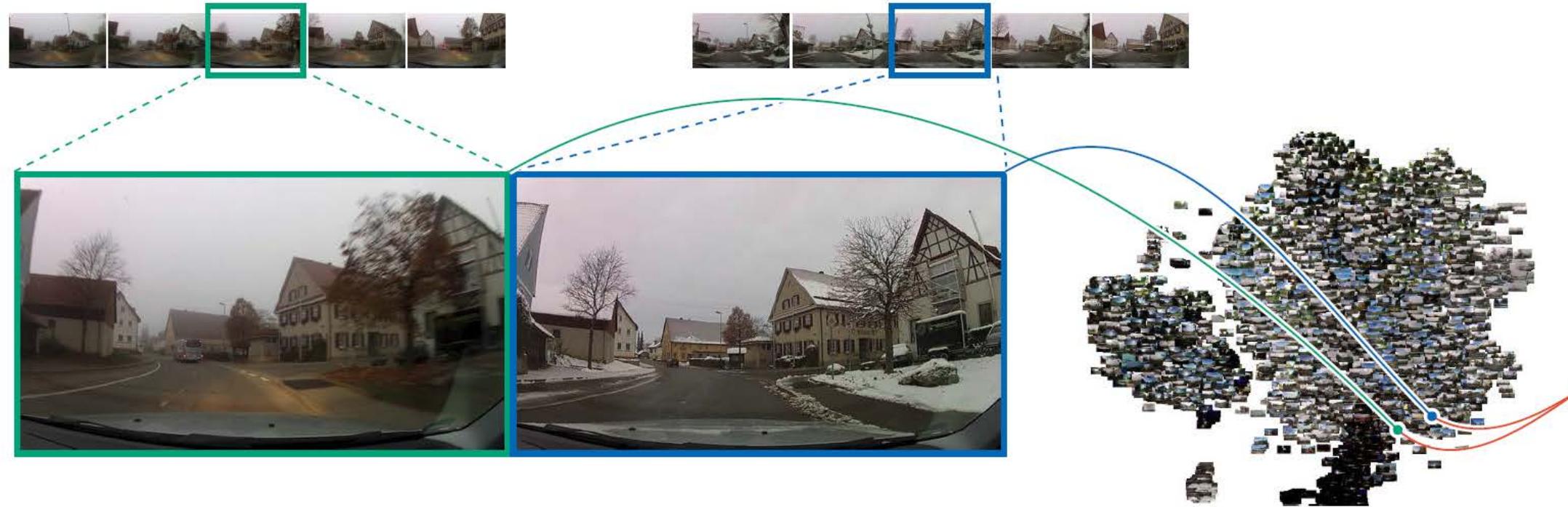
- Assume ground-truth annotations are given
- Optimize triplet loss



$$\max\{0, m + \|\Phi_a - \Phi_p\|_2^2 - \|\Phi_a - \Phi_n\|_2^2\}$$



## Encoding each Frame

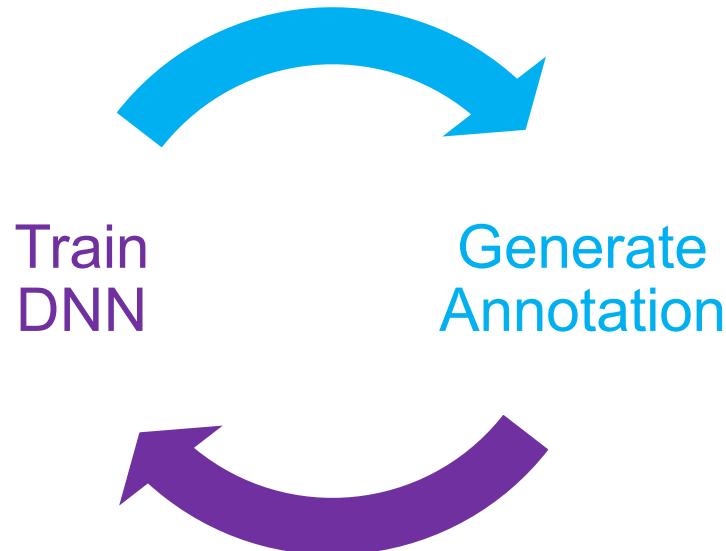


- Encoding clearly is not perfect
- Interplay of CNN encoding and shortest path on cleaned matrix render it a stable process



# Data-Driven Approaches without Annotations

- Chicken-Egg-Problem

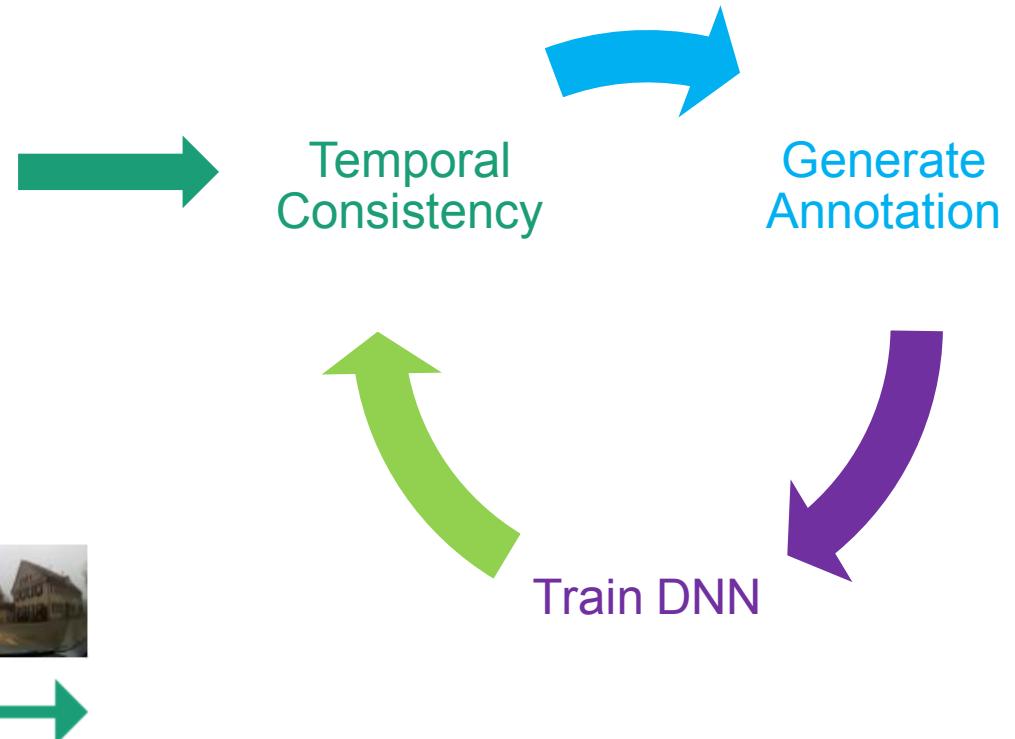




# Data-Driven Approaches without Annotations

- Solve the Chicken-Egg-Problem

- Causality of Time

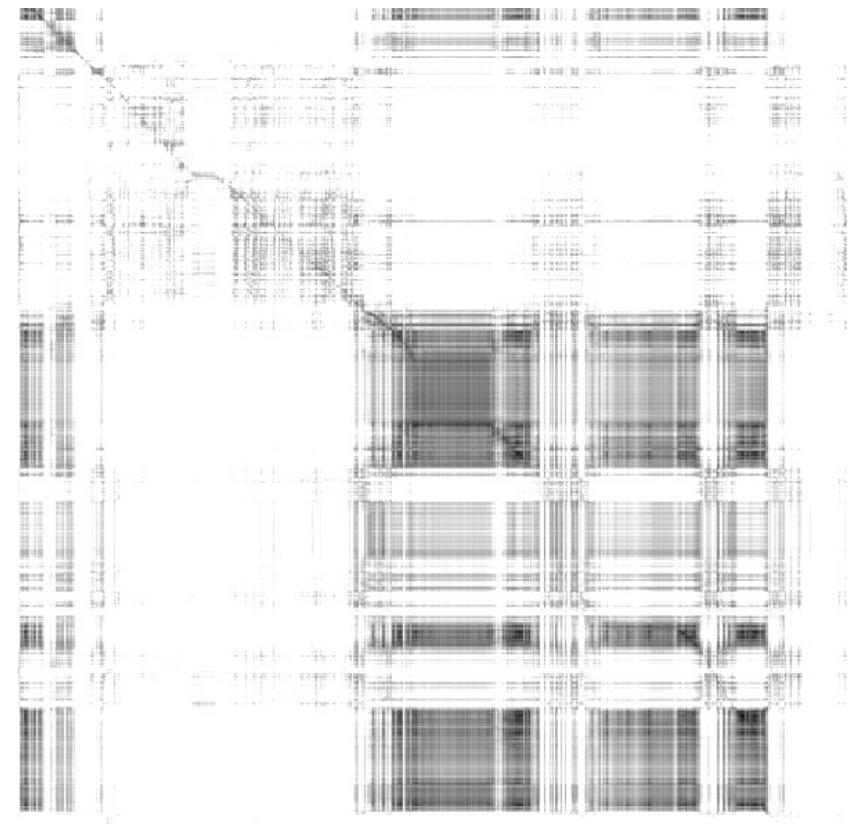




# Curriculum Learning

- Start with examples which are not too difficult and not too obvious:

1. Intra-Video Sampling  
consecutive frames are similar
2. Inter-Video Sampling  
ensure temporal consistency to reject false positives
3. Transitive Inter-Video Sampling  
 $\nu/\sim$  is equivalence relation  $a \sim b \sim c$







# Important Features

- Why are local features not useful?
- What is the network looking for?

HARRIS



used in Evangelidis et al. 2011, 2013

SIFT



used in Wang et al. 2014



# Learned Features (Saliency Maps)

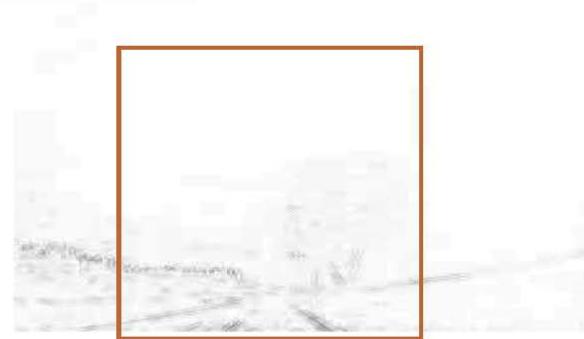
Input



ResNet (ImageNet)

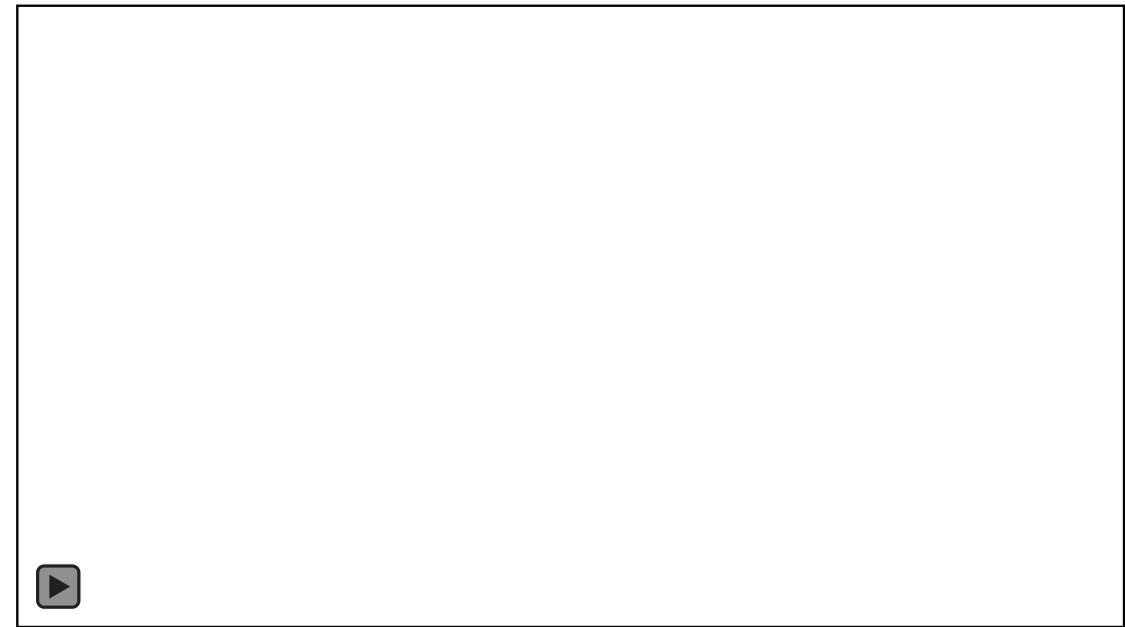


Ours





# Saliency Maps









# Robustness

**after tree-felling**



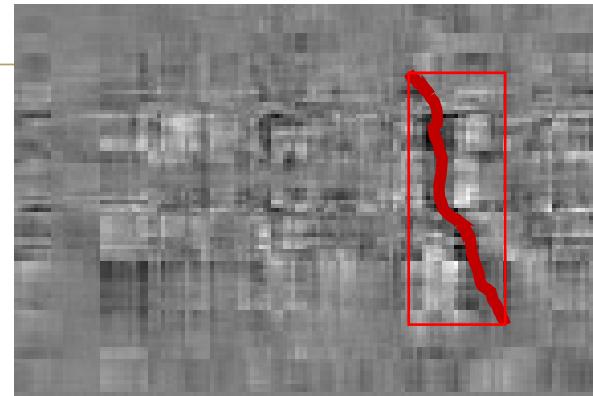
**before tree-felling**





## Application to other Events

- Golf swing – same neural network
- Event boundaries visible





# Global Illumination

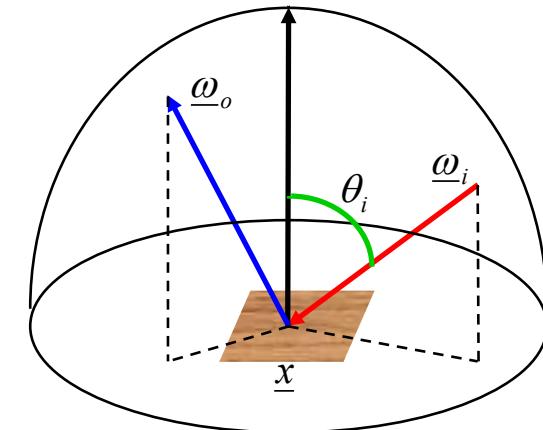


# Surface Radiance

$$\begin{aligned} L(\underline{x}, \underline{\omega}_o) &= L_e(\underline{x}, \underline{\omega}_o) + L_r(\underline{x}, \underline{\omega}_o) \\ &= L_e(\underline{x}, \underline{\omega}_o) + \int_{\Omega_+} f_r(\underline{\omega}_i, \underline{x}, \underline{\omega}_o) L(\underline{x}, \underline{\omega}_i) \cos \theta_i d\underline{\omega}_i \end{aligned}$$

- Visible surface radiance
  - Surface position
  - Outgoing direction
  - Incoming illumination direction
- Self-emission
- Reflected light
  - Incoming radiance from all directions
  - Direction-dependent reflectance  
(BRDF: bidirectional reflectance distribution function)

$$\begin{array}{l} L(\underline{x}, \underline{\omega}_o) \\ \underline{x} \\ \underline{\omega}_o \\ \underline{\omega}_i \\ L_e(\underline{x}, \underline{\omega}_o) \end{array}$$



$$\begin{array}{l} L(\underline{x}, \underline{\omega}_i) \\ f_r(\underline{\omega}_i, \underline{x}, \underline{\omega}_o) \end{array}$$

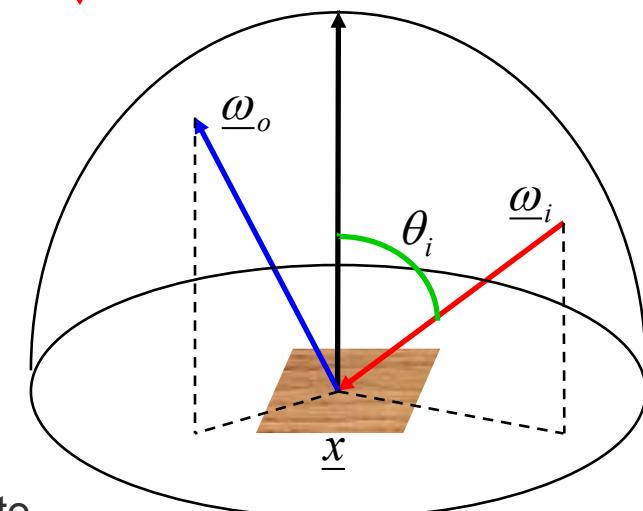


# (Surface) Rendering Equation

- In Physics: Radiative Transport Equation
- Expresses energy equilibrium in scene

$$L(\underline{x}, \underline{\omega}_o) = L_e(\underline{x}, \underline{\omega}_o) + \int_{\Omega} f_r(\underline{\omega}_i, \underline{x}, \underline{\omega}_o) L(\underline{x}, \underline{\omega}_i) \cos \theta_i d\underline{\omega}_i$$

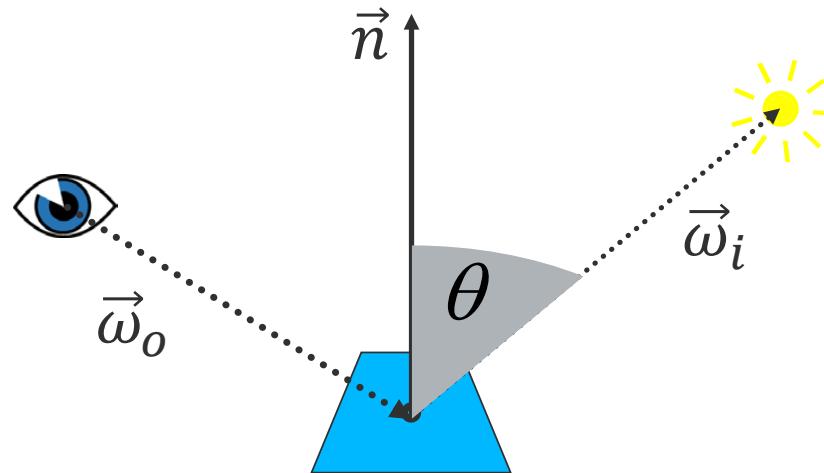
- total radiance = emitted radiance + reflected radiance
- First term: emissivity of the surface
  - non-zero only for light sources
- Second term: reflected radiance
  - integral over all possible incoming directions of irradiance times angle-dependent surface reflection function
- Fredholm integral equation of 2nd kind
  - unknown radiance appears on lhs and inside the integral
  - Numerical methods necessary to compute approximate solution





# Sampling the Incoming Light

- Consider area foreshortening (cosine)

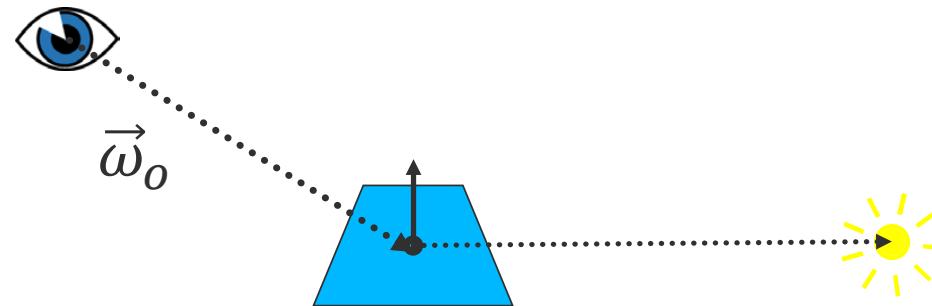


$$L_r = \int_{\Omega} f_r(\vec{\omega}_o, \vec{\omega}_i) L_i \cos(\theta) d\omega_i = \int_{\Omega} f_r(\vec{\omega}_o, \vec{\omega}_i) L_i \langle \vec{n}, \vec{\omega}_i \rangle d\omega_i$$



# Sampling the Incoming Light

- Continuous integral

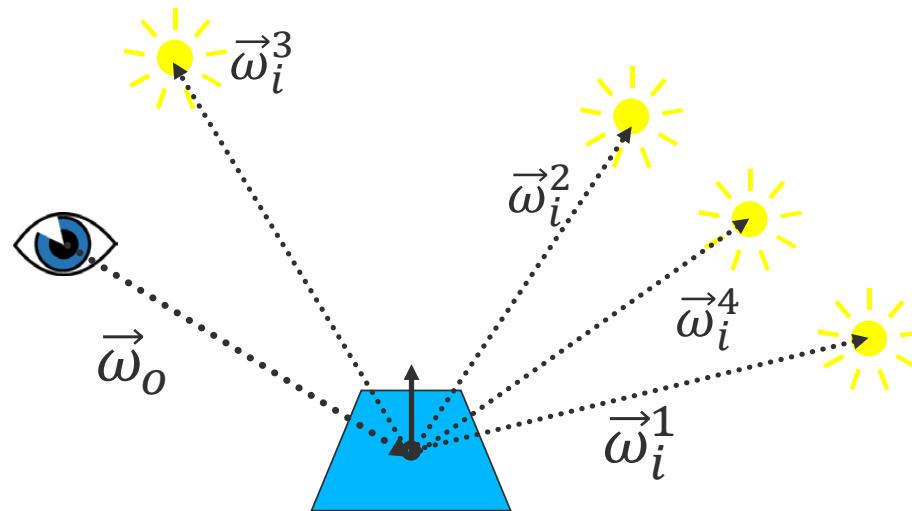


$$L_r = \int_{\Omega} f_r(\vec{\omega}_o, \vec{\omega}_i) L_i \cos(\theta) d\omega_i$$



# Sampling Light Sources

- Consider discrete samples of the path space
- Multiple Light Sources:
  - add their contributions



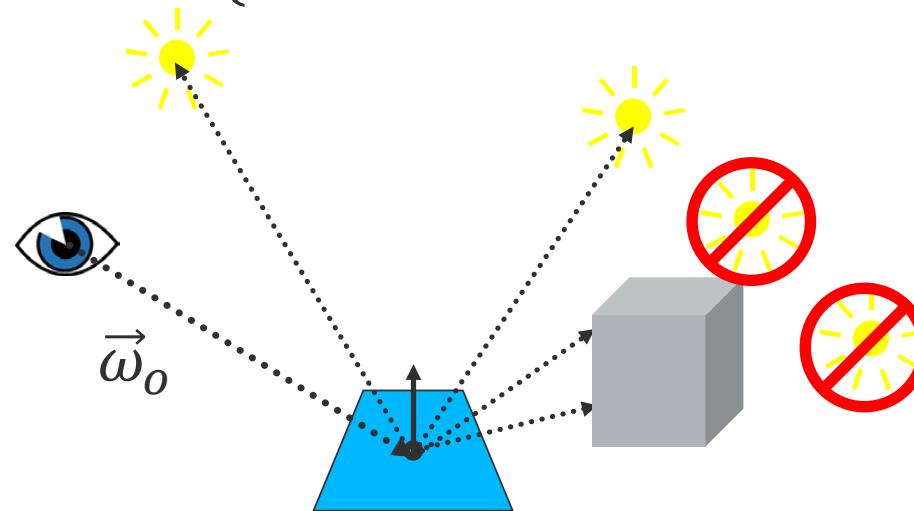
$$L_r = \sum_N f_r(\vec{\omega}_o, \vec{\omega}_i) L_i^k \cos(\theta^k) \approx \int_{\Omega} f_r(\vec{\omega}_o, \vec{\omega}_i) L_i \cos(\theta) d\omega_i$$



# Sampling Light Sources w/ Occlusions

- Visibility function

$$v(\vec{\omega}_i) = \begin{cases} 0 & \text{source is occluded} \\ 1 & \text{source is visible} \end{cases}$$



$$L_r = \sum_N f_r(\vec{\omega}_o, \vec{\omega}_i) v(\vec{\omega}_i) L_i^k \cos(\theta^k)$$



# Approximation – Ambient Occlusion

- Approximation for global illumination effects
- Only consider visibility in a uniformly lit room
- AO = fraction of the hemisphere that is not occluded

$$L_r = \int_{\Omega} f_r(\vec{\omega}_o, \vec{\omega}_i) v(\vec{\omega}_i) L_i \cos(\theta) d\omega_i$$

$$L_r \approx L_{AO} = \rho \cdot AO \cdot \int_{\Omega} L_i \cos(\theta) d\omega_i = \rho \cdot AO \cdot E$$

$$AO = \frac{1}{2\pi} \int_{\Omega} v'(\vec{\omega}_i) d\omega_i$$

$v'(\vec{\omega}_i) = v(\vec{\omega}_i)$   
up to some  
max. distance



## Example: Lighting



$$L_r \approx L_{AO} = \rho \cdot AO \cdot \int_{\Omega} L_i \cos(\theta) d\omega_i = \rho \cdot AO \cdot E$$



## Example: Lighting & AO



$$L_r \approx L_{AO} = \rho \cdot AO \cdot \int_{\Omega} L_i \cos(\theta) d\omega_i = \rho \cdot AO \cdot E$$



## Example: Lighting & Shading & AO



$$L_r \approx L_{AO} = \rho \cdot AO \cdot \int_{\Omega} L_i \cos(\theta) d\omega_i = \rho \cdot AO \cdot E$$



## Example: Lighting & Shading

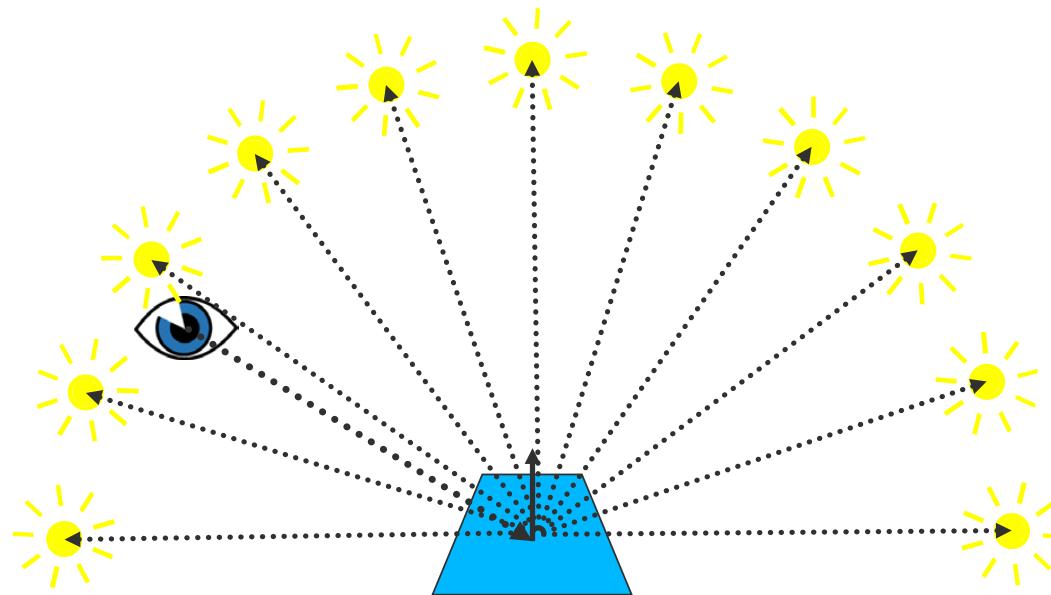


$$L_r \approx L_{AO} = \rho \cdot AO \cdot \int_{\Omega} L_i \cos(\theta) d\omega_i = \rho \cdot AO \cdot E$$



# Ambient Occlusion

- Only consider visibility in a uniformly lit room
- AO = fraction of the hemisphere that is not occluded

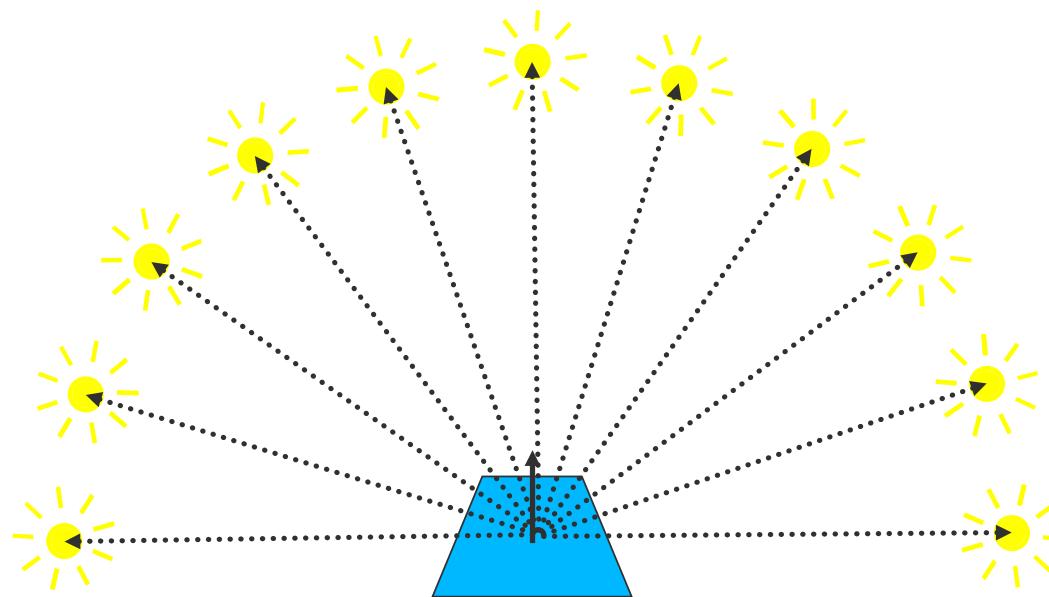


$$AO = \frac{1}{N} \sum_N v(\vec{\omega}_i^k)$$



# Ambient Occlusion

- Fraction of the hemisphere that is not occluded
- Independent on viewing direction

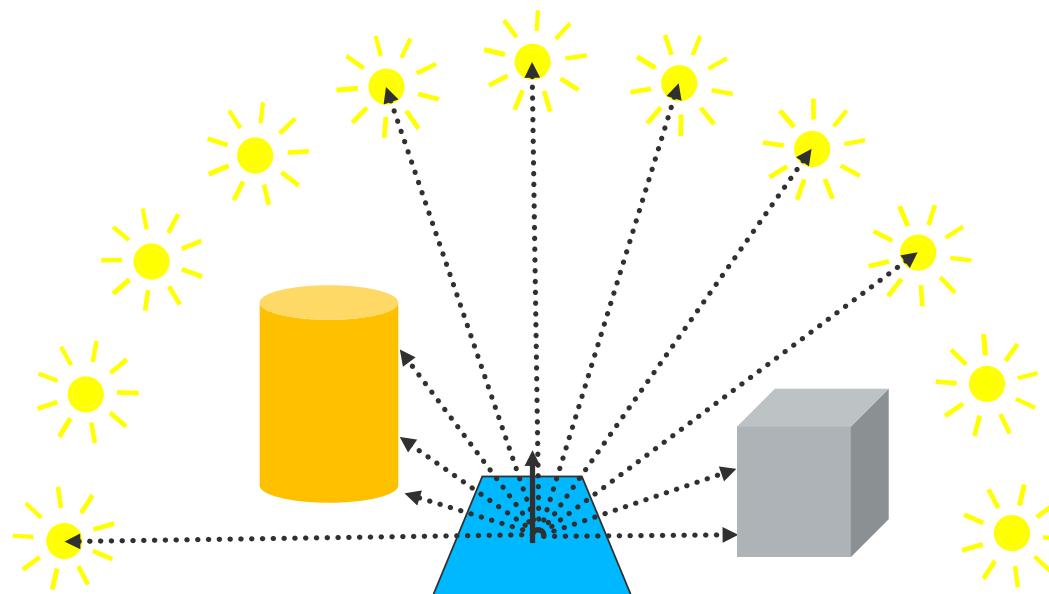


$$AO = \frac{1}{N} \sum_N v(\vec{\omega}_i^k) = 1 \quad \text{completely unoccluded}$$



# Ambient Occlusion

- Fraction of the hemisphere that is not occluded

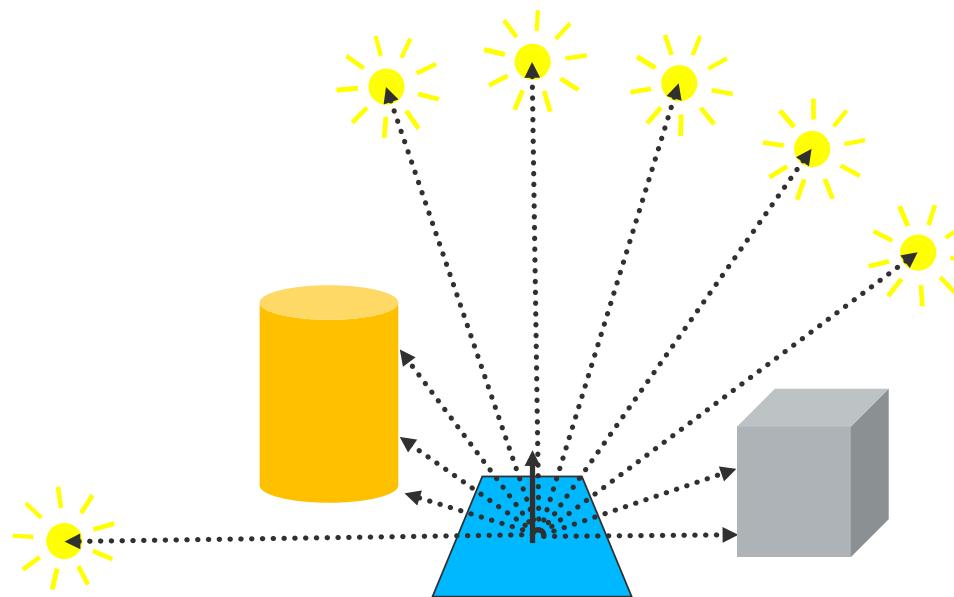


$$AO = \frac{1}{N} \sum_N v(\vec{\omega}_i^k)$$



# Ambient Occlusion

- Fraction of the hemisphere that is not occluded

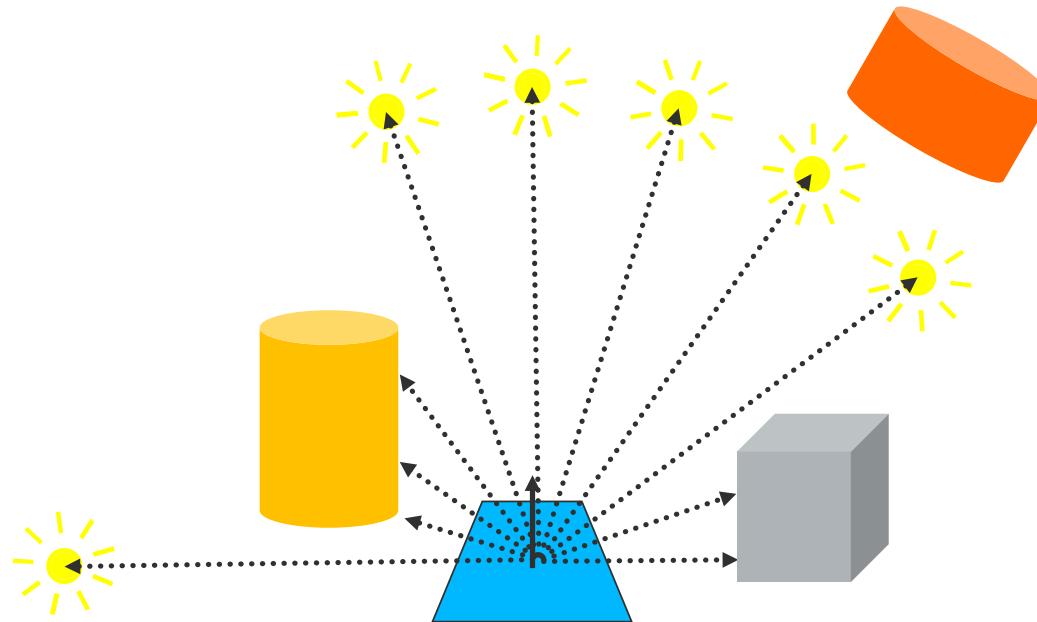


$$AO = \frac{1}{N} \sum_N v(\vec{\omega}_i^k)$$



# Ambient Occlusion

- Fraction of the hemisphere that is not occluded
- Only consider close by objects



$$AO = \frac{1}{N} \sum_N v(\vec{\omega}_i^k)$$



# Ambient Occlusion

- Almost looks like global light transport





## The problem

- Variance appears as noise at low sample rates





# General Framework for Monte Carlo Filtering



*scene by Herminio Nieves and tf3dm.com user ysup12*



# Observations

---

- After only a few minutes, we have a lot of information about the scene...



**scene rendered in a few minutes (8 samples/pixel)**

*scene by M. Leal Llaguno*



# Observations

- Can we filter this to get a result as good as a reference image?



scene rendered in over 24 hrs (8,192 samples/pixel)

scene by M. Leal Llaguno

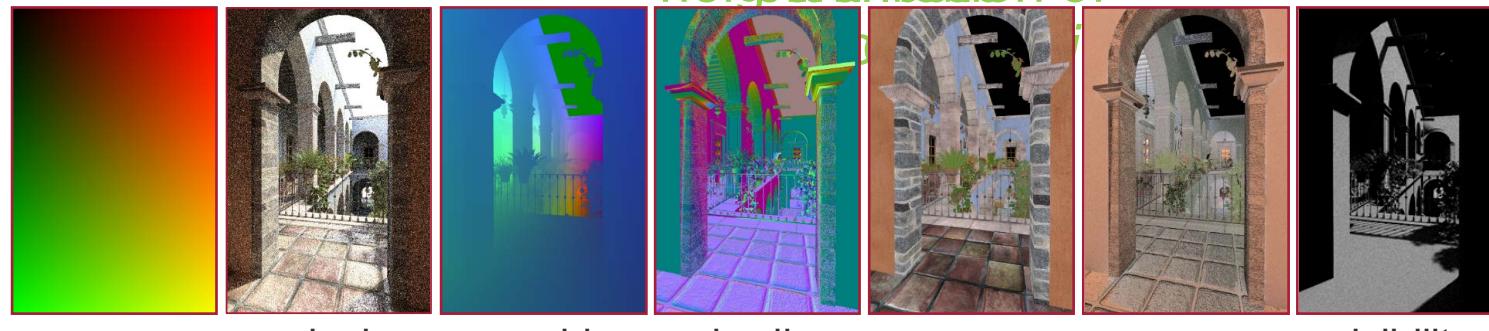


# Filter

$$\hat{\mathbf{c}}_i = \frac{\sum_{j \in \mathcal{N}(i)} d_{i,j} \bar{\mathbf{c}}_j}{\sum_{j \in \mathcal{N}(i)} d_{i,j}}$$



neighborhood of



screen  
position

pixel  
color

world  
positions

shading  
normals

texture  
1<sup>st</sup> bounce

texture  
2<sup>nd</sup> bounce

visibility



# Cross-bilateral filter

$$\hat{\mathbf{c}}_i = \frac{\sum_{j \in \mathcal{N}(i)} d_{i,j} \bar{\mathbf{c}}_j}{\sum_{j \in \mathcal{N}(i)} d_{i,j}}$$

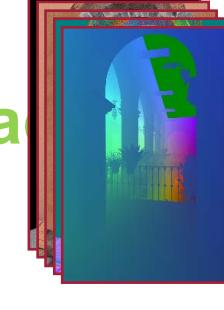
$$d_{i,j} = \exp \left[ - \frac{\|\bar{\mathbf{p}}_i - \bar{\mathbf{p}}_j\|^2}{2\alpha_i^2} \right] \quad \beta_i^2 \quad \gamma_{k,i}^2$$

- $\alpha_i^2, \beta_i^2$  term  $\gamma_{k,i}^2$  are variance of each parameter
- 

screen position



pixel color



additional features



# Why per pixel parameters?



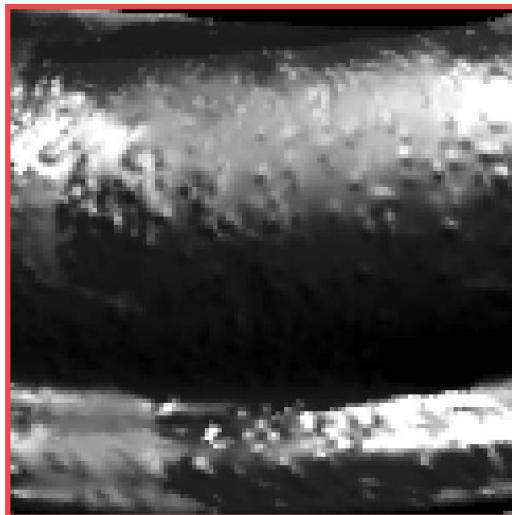
input MC 8 samples/pixel



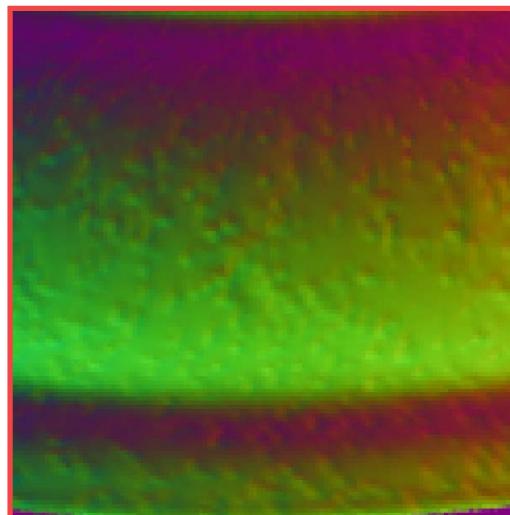
shading normals



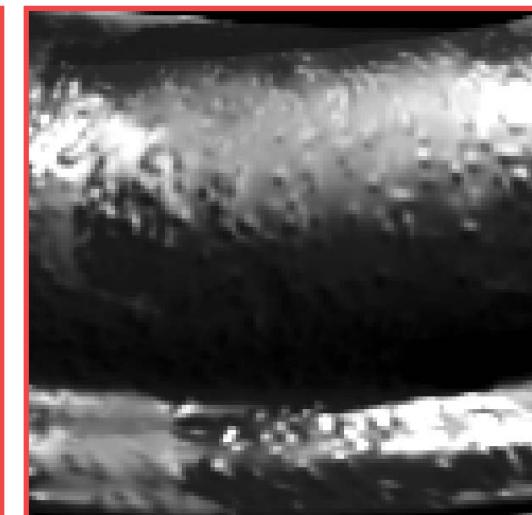
# Why per pixel parameters?



input MC  
8 samples/pixel



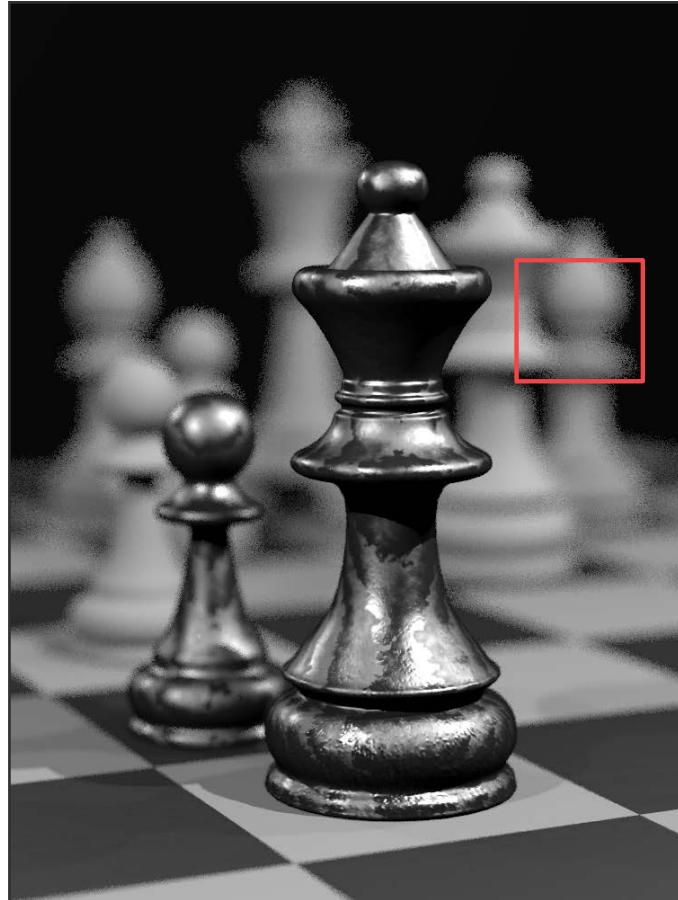
shading normals



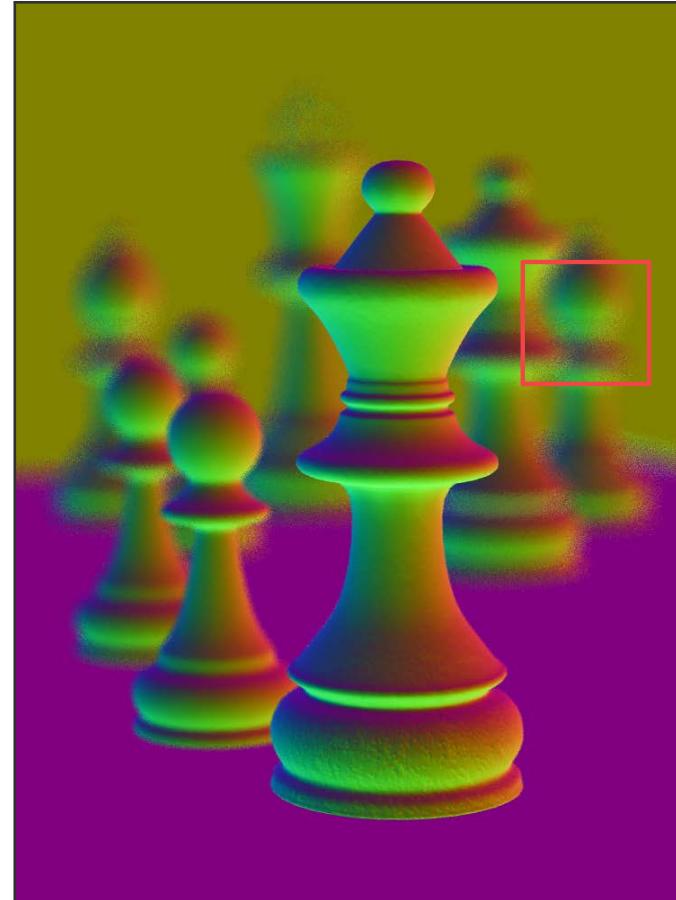
ground truth



# Why per pixel parameters?



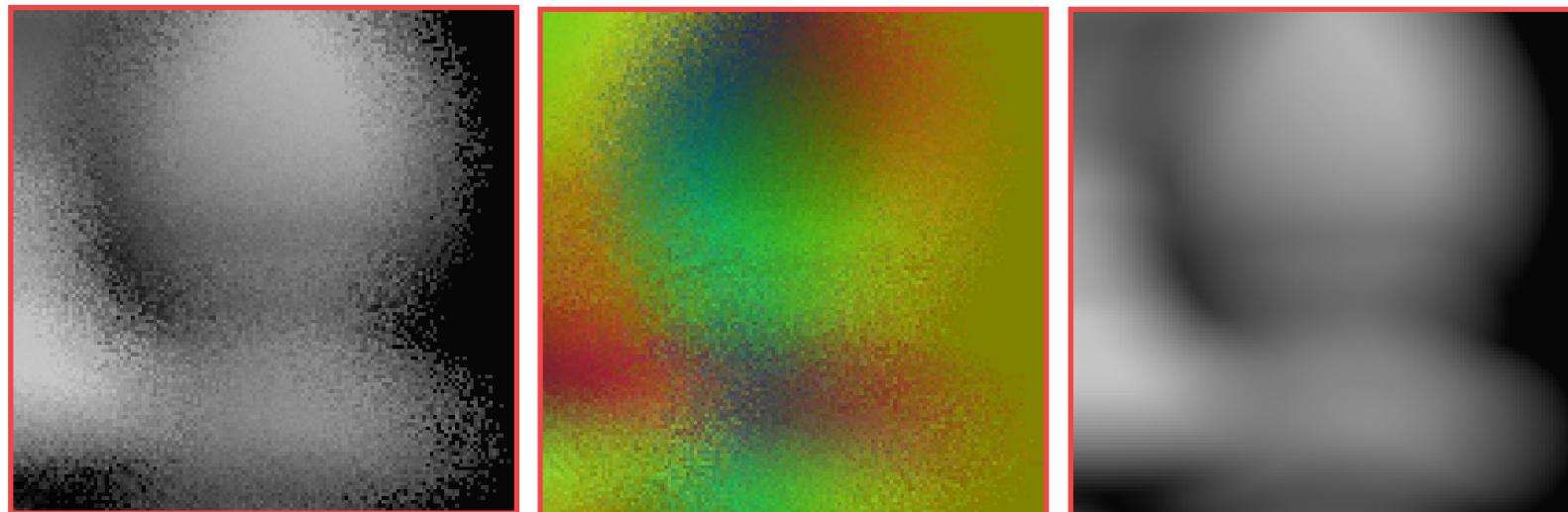
input MC 8 samples/pixel



shading normals



# Why per pixel parameters?



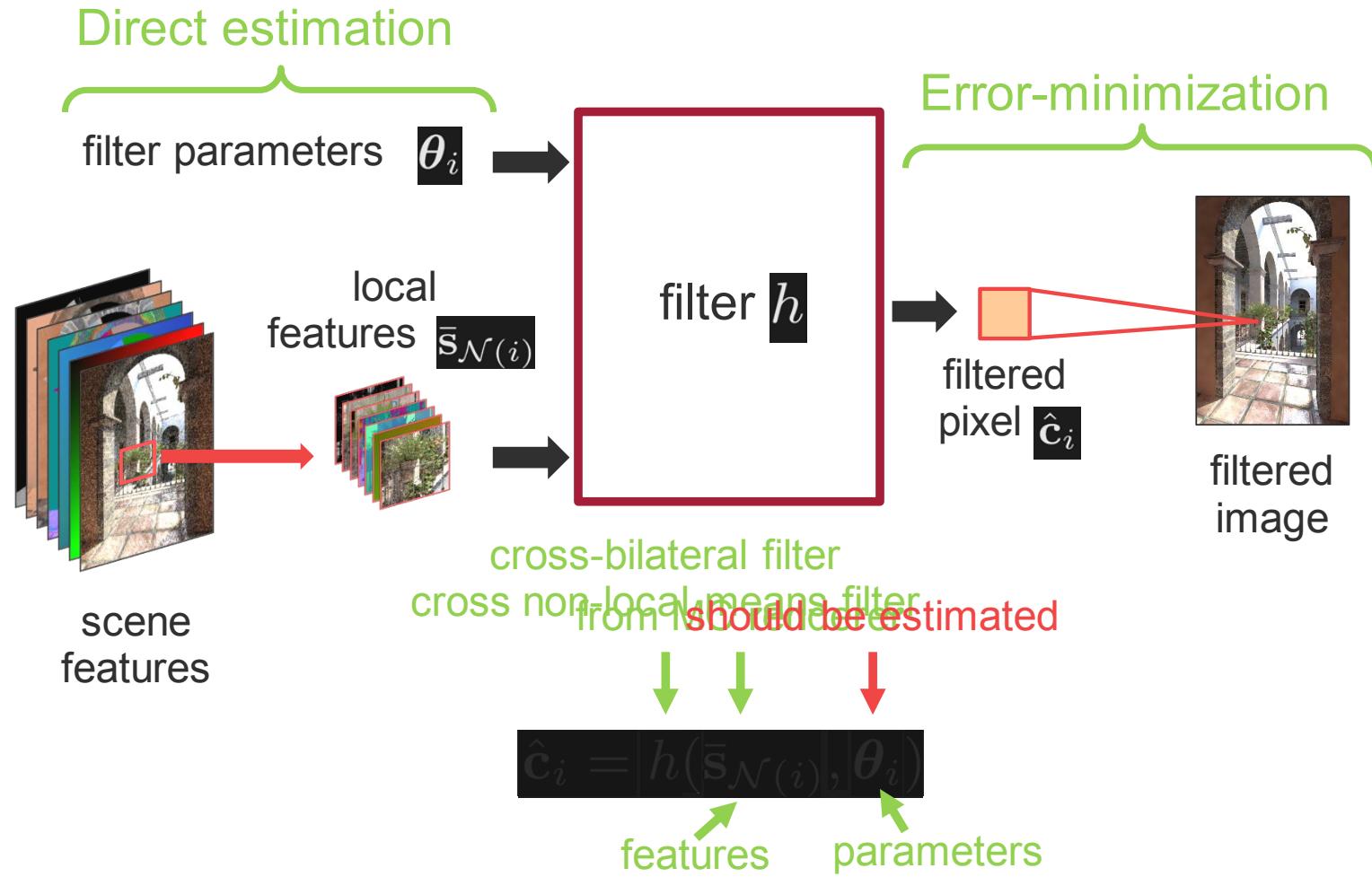
input MC  
8 samples/pixel

shading normals

ground truth

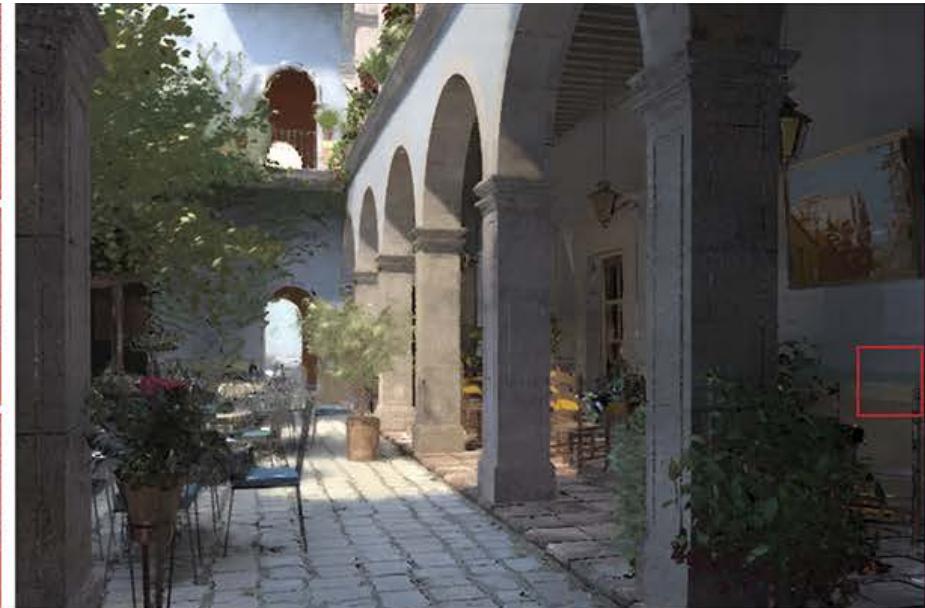


# Problem formulation

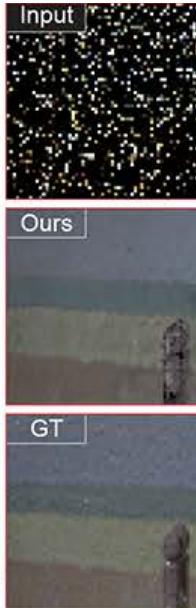




Our result with a cross-bilateral filter (4 spp)



Our result with a non-local means filter (4 spp)



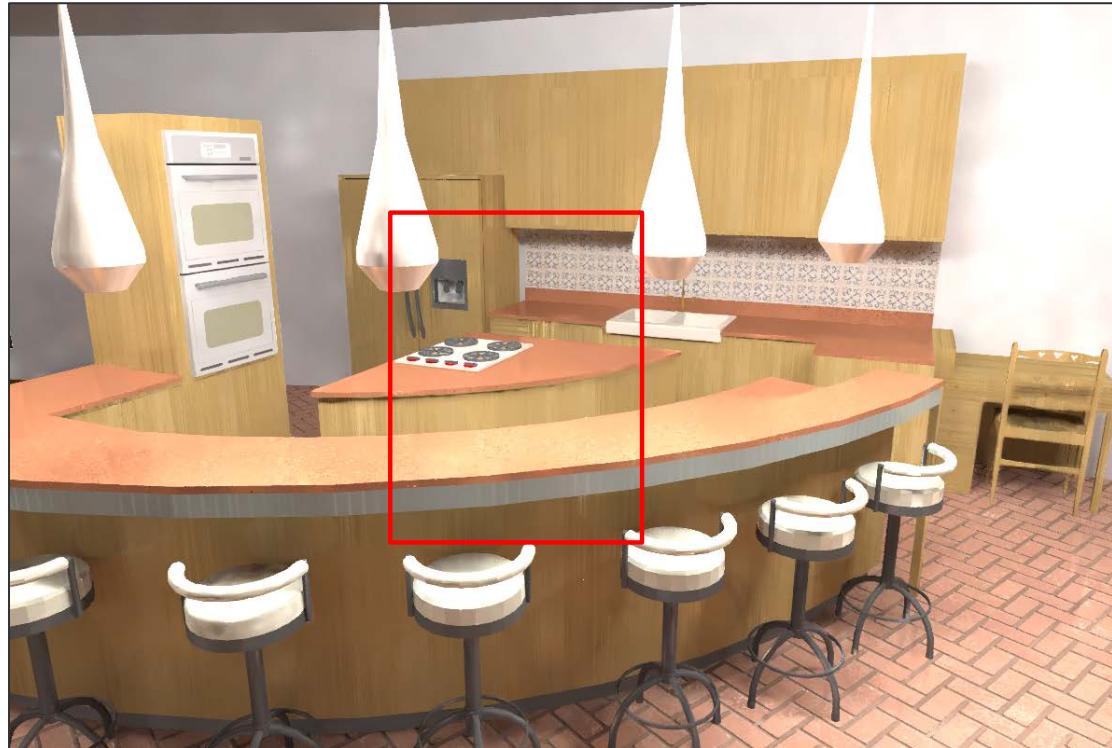
# A Machine Learning Approach for Filtering Monte Carlo Noise

Kalantari et al. SIGGRAPH 2015



# Sample result

Path traced scene



4 samples/pixel  
(48.9 sec)

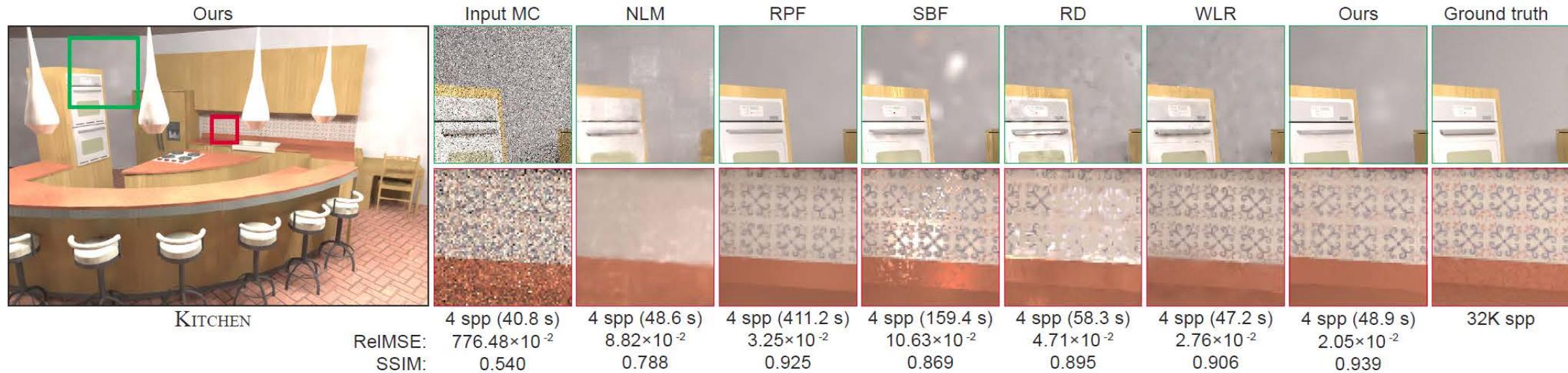
[Kalantari et al. 2015]



using only post-process filter!



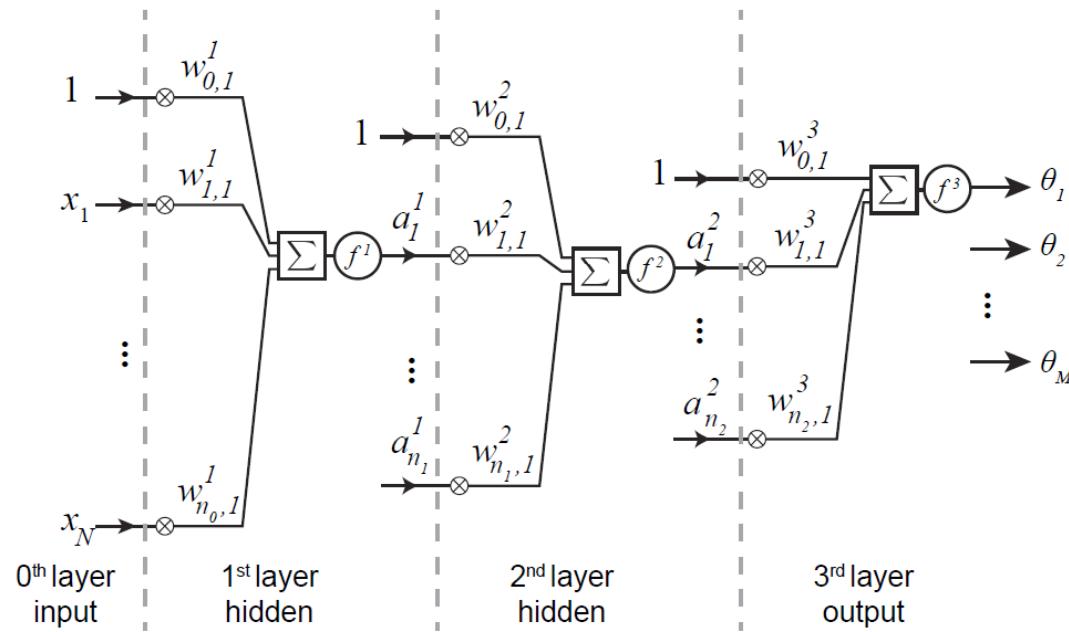
# Results





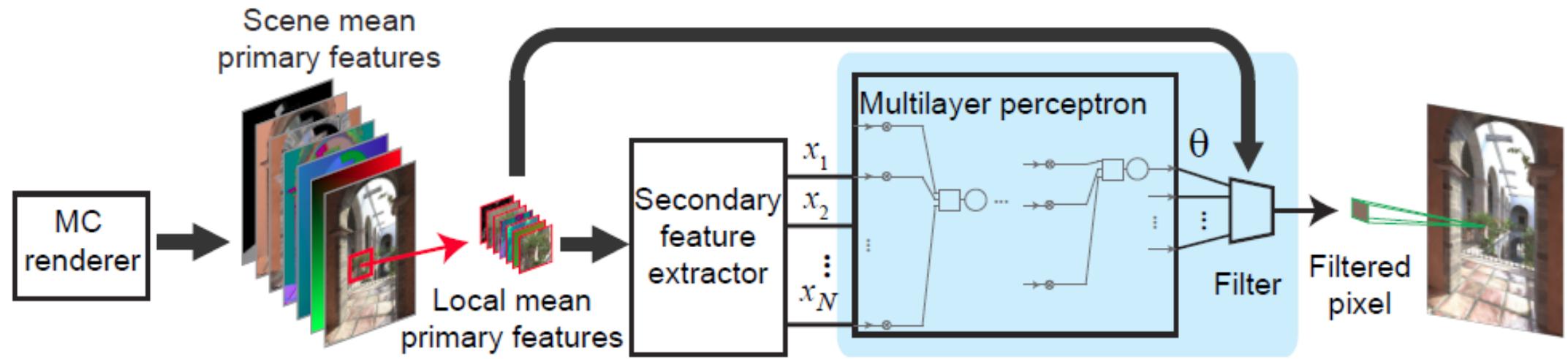
# MLP for Denoising

- Simple 3-layer MLP





# Incorporating primary Features







# Denoising with Kernel Prediction and Asymmetric Loss Functions

Vogels et al. SIGGRAPH 2018

# Kernel-predicting Convolutional Network (KPCN)

- training to minimize the average distance between the ground-truth  $r$  and the denoised image  $d = g(x; \theta)$  via a loss function  $l$

$$\hat{\theta} = \operatorname{argmin}_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(r^n, g(x^n; \theta))$$

- symmetric mean absolute percentage error

$$\ell(\mathbf{r}, \mathbf{d}) = \frac{1}{3|\mathcal{I}|} \sum_{p \in \mathcal{I}} \sum_{c \in C} \frac{|\mathbf{d}_{p,c} - \mathbf{r}_{p,c}|}{|\mathbf{d}_{p,c}| + |\mathbf{r}_{p,c}| + \varepsilon}$$

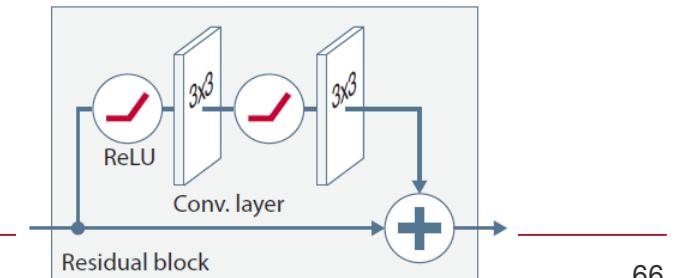
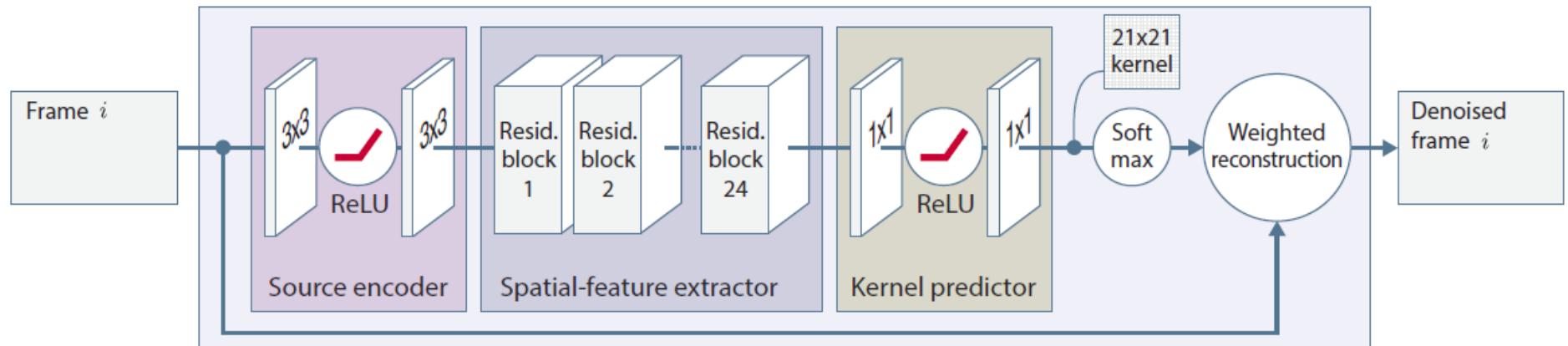


# Kernel-predicting Convolutional Network (KPCN)

- Network predicts 21x21 filter kernel for each neighborhood followed by softmax.

$$d_p = g_p(x; \theta) = \sum_{q \in N(p)} w_{pq} c_q$$

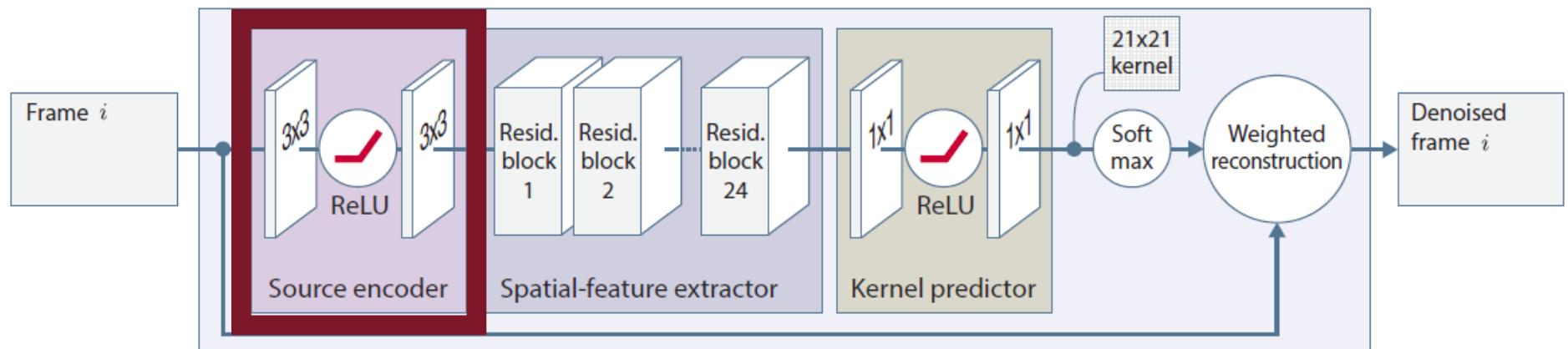
$$w_{pq} = \frac{\exp([z_p]_q)}{\sum_{q' \in N(p)} \exp([z_p]_{q'})}$$





# Source-aware Encoder

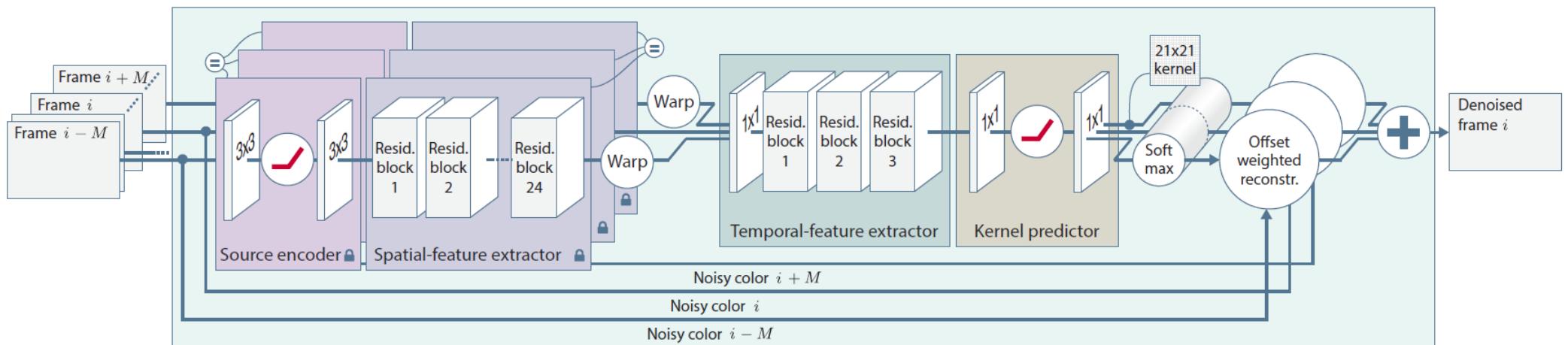
- Simple adaptation to different renderer providing different input modalities
- Can incorporate different auxilliary buffers (color, depth, geometry, normal, texture, ....)
- Simple to retrain



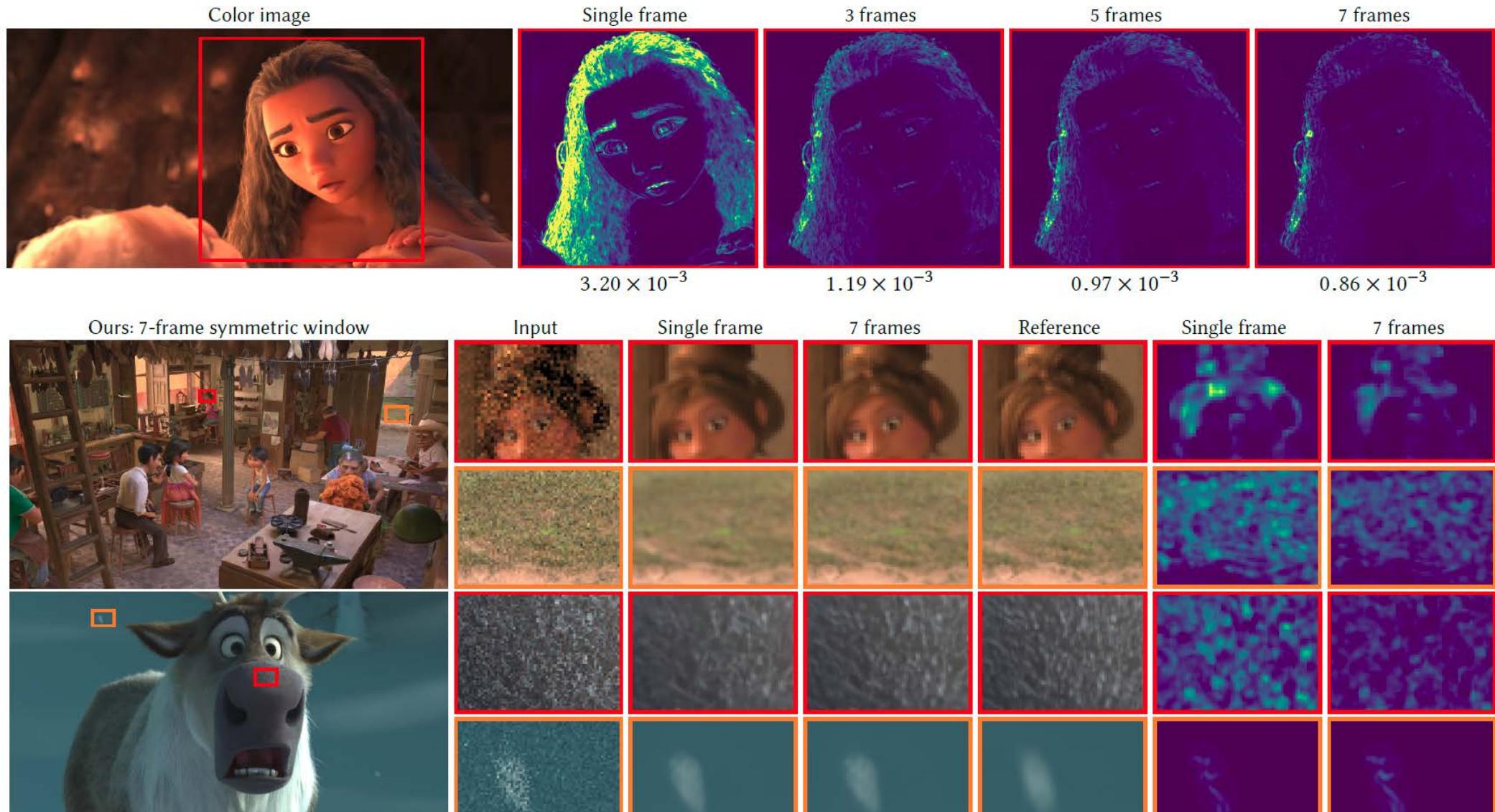


# Temporal denosing

- Jointly predict kernels for eachs consecutive frame
- Warping based on known geometry (we are denosing images produced by a renderer) or by optical flow
- Train source encoder and feature extractor for individual frames first
- Then train temporal component



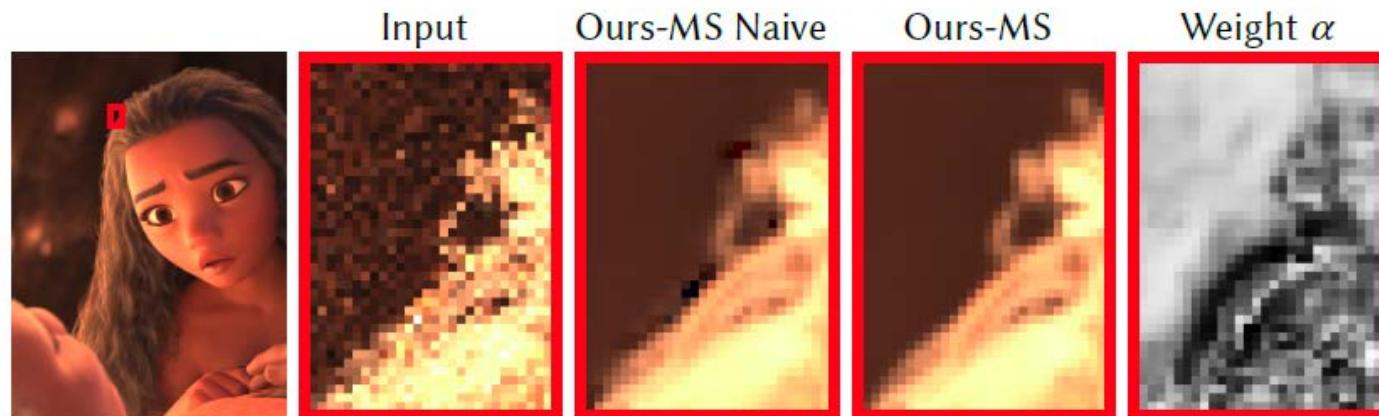
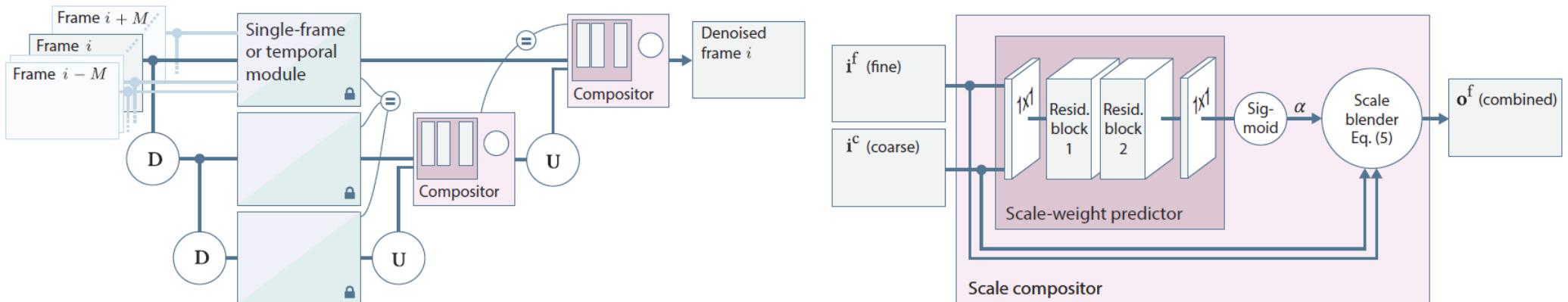
# Temporal Window Size





# Multi-scale Prediction

- Single scale denoising might suffer from low-frequency noise
- Compose final result by weighted scale blending





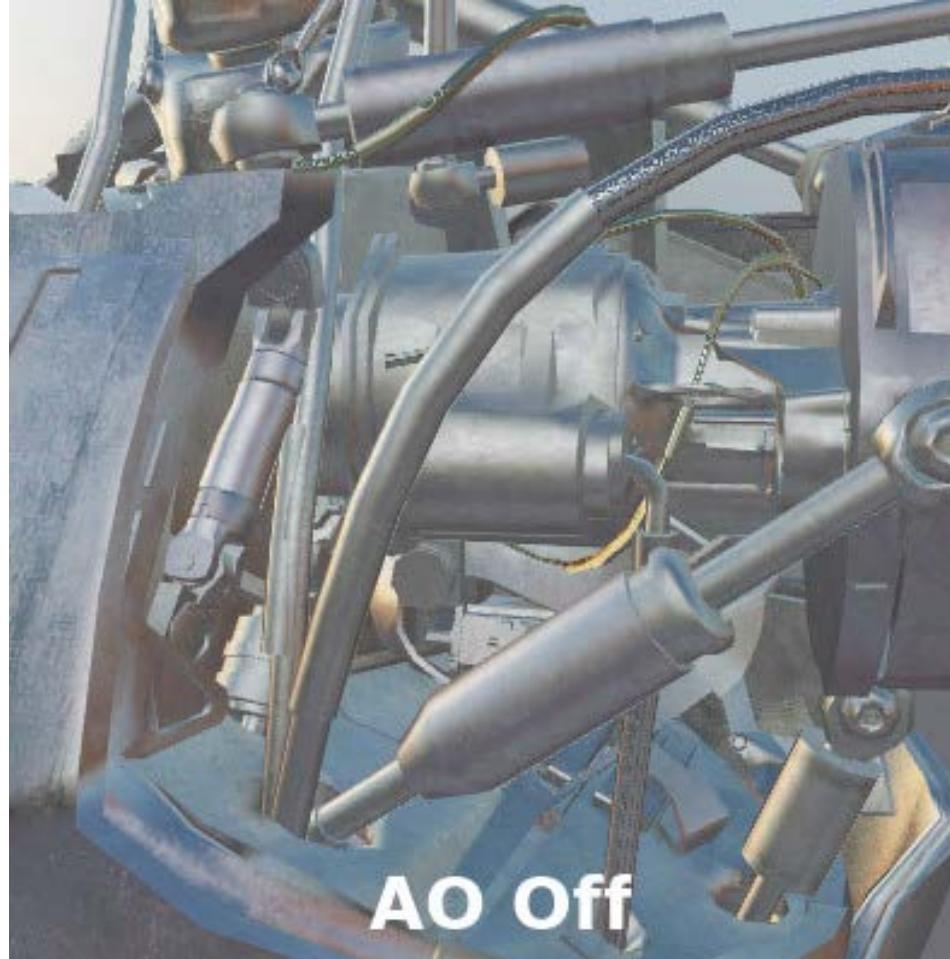


# Neural Network Ambient Occlusion

Holden et al. SIGGRAPH Asia 2016



# AO approximates Global Illumination





# AO approximates Global Illumination



## Cost of AO evaluation

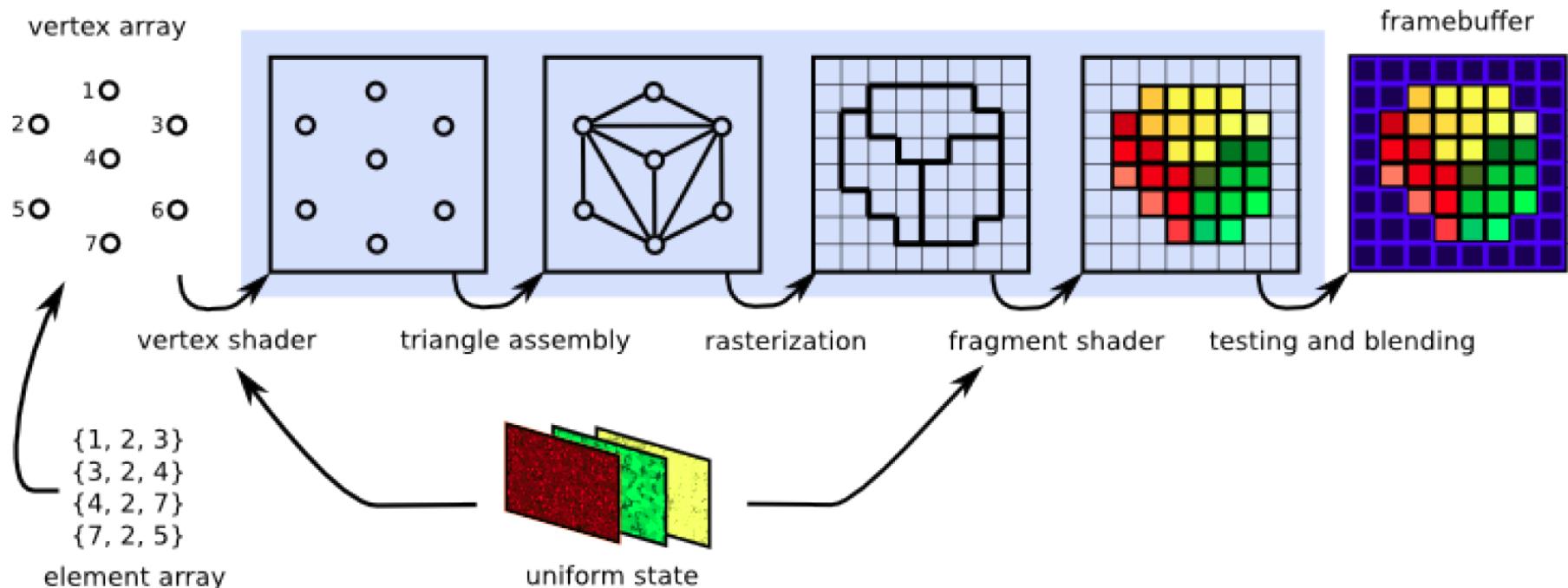
---

- Raytracing AO still slow
- (256x256 px, 1 sample per pixel, 1024 samples per point = 67 108 864 rays) 2.5 seconds
- Too slow for games
- Idea: Utilize information that is already available from the G-Buffer in the Deferred Shading pipeline



# G-Buffer

- Directly supported by GPUs

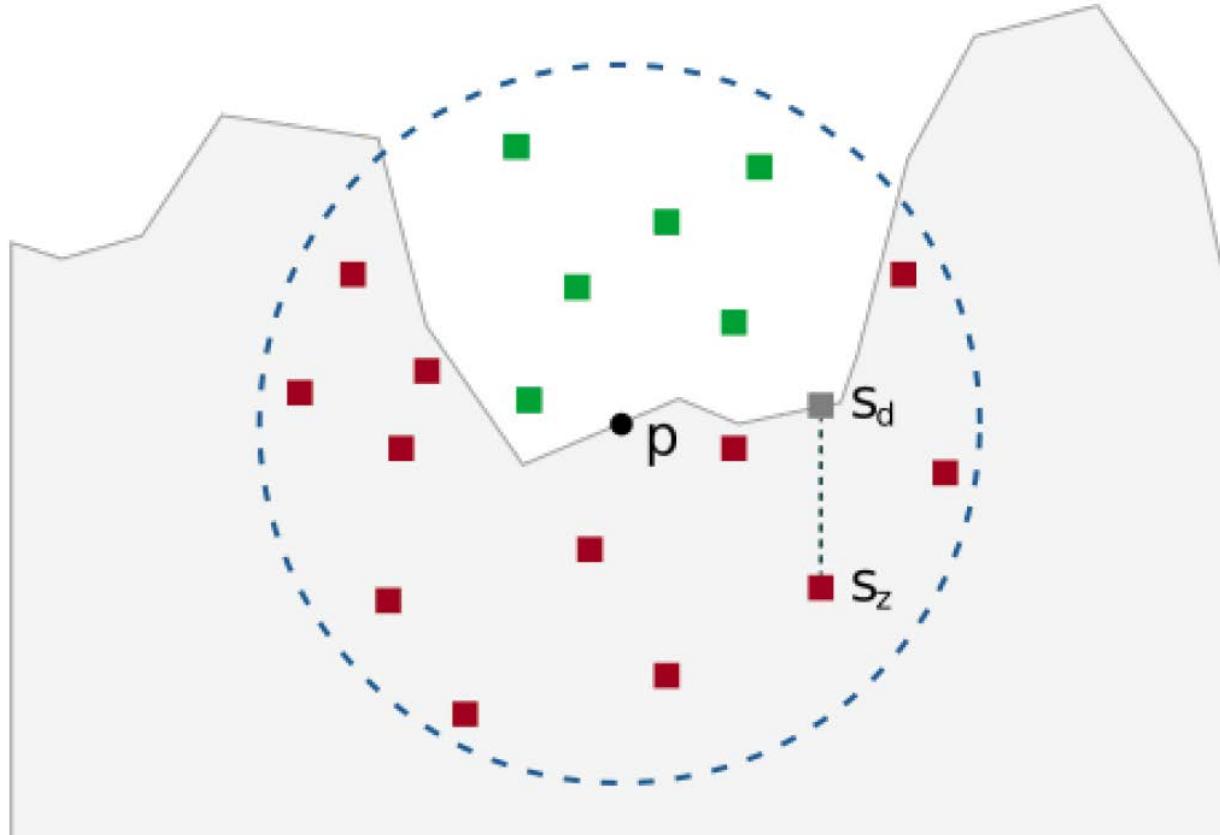


Source: <https://www.enlightenment.org/playground/evas-gl.md>



# Screen-space Ambient Occlusion

- Idea: Percentage of samples placed behind the depth map.

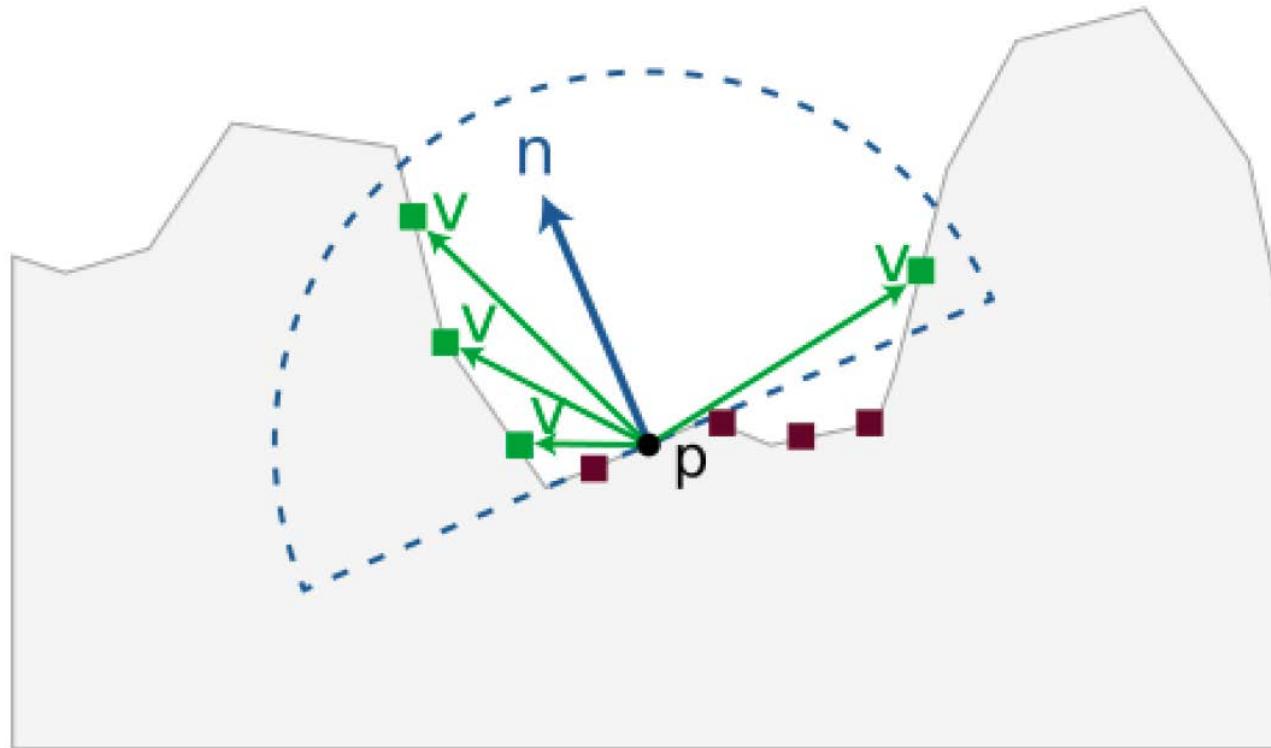


Source: F. Aalund and J.A. Bärentzen. A comparative study of screen-space ambient occlusion methods. Bachelor Thesis, Technical University of Denmark Informatics and Mathematical Modelling, 2013



# Alchemy Ambient Occlusion

- Idea: Sum of weighted sample point distances.

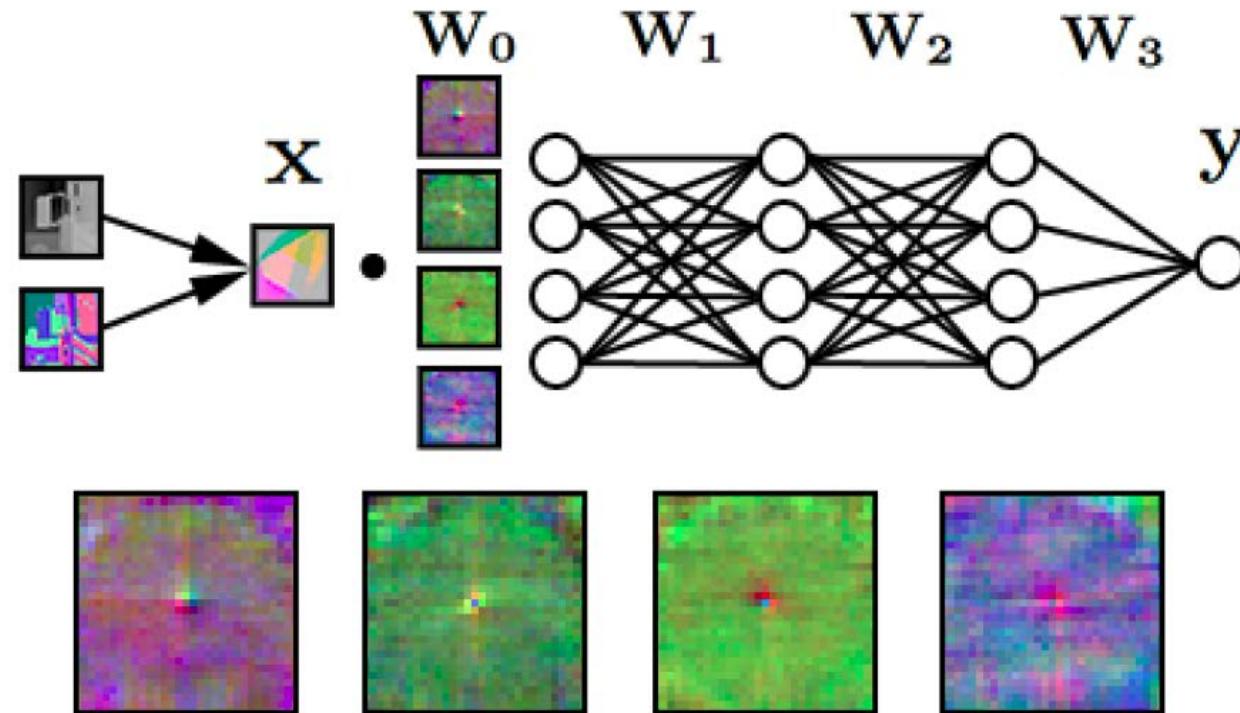


Source: F. Aalund and J.A. Bærentzen. A comparative study of screen-space ambient occlusion methods. Bachelor Thesis, Technical University of Denmark Informatics and Mathematical Modelling, 2013



# NNAO: Simple MLP for Ambient Occlusion Prediction

- Predict the AO value of the central pixel of one region
- Input: prepared neighborhood from screen-space buffers
  - Absolute values and difference to central pixel of depth and normal

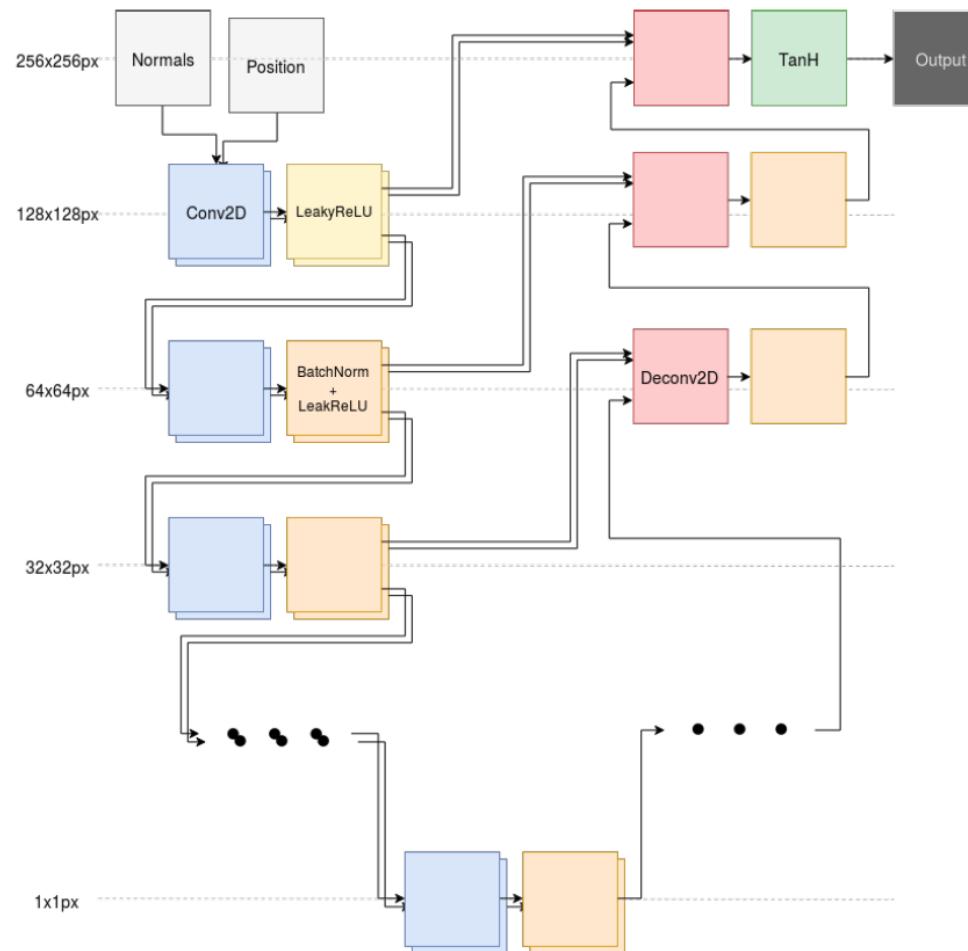




- 
- Video



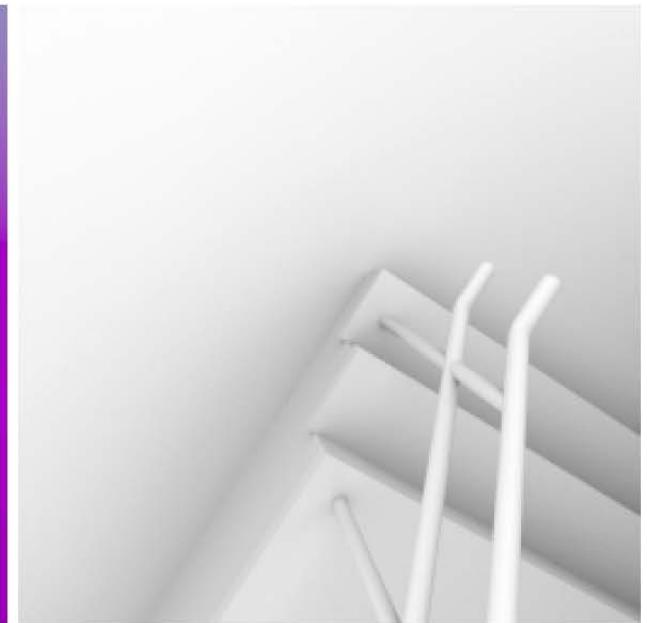
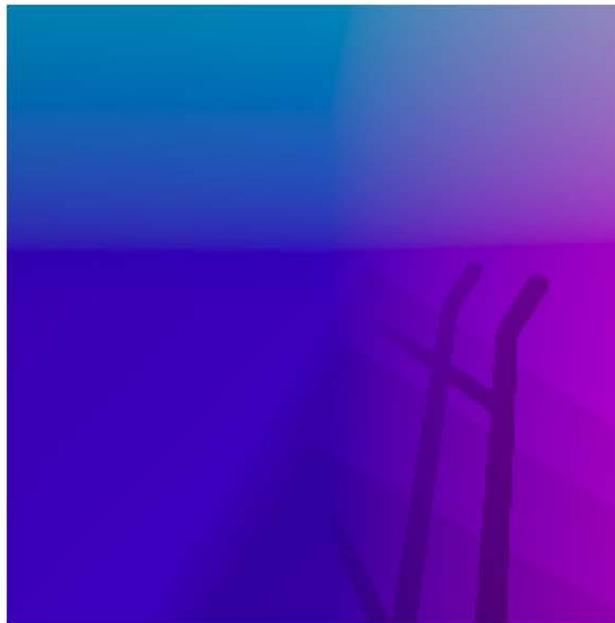
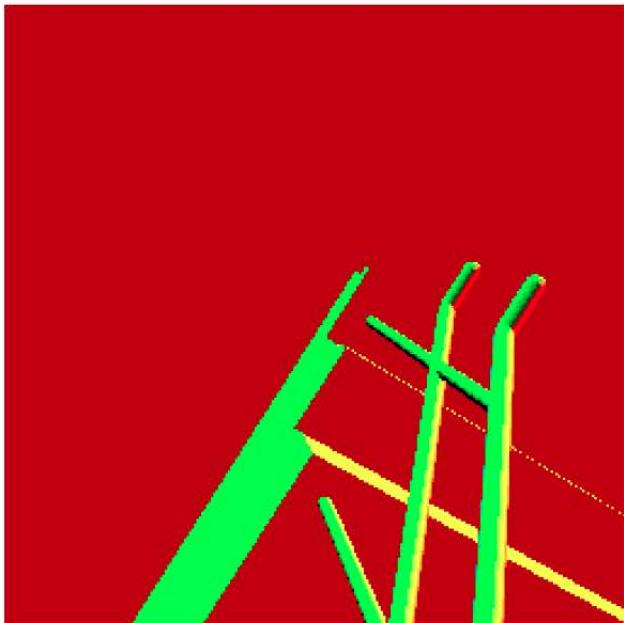
# Reimplementation (Jens Schindel)





# Training Data

- Training Set 358650 image pairs
- Validation Set 72000 image pairs (20%)
- Test Set 27002 image pairs (7.5%)



[Jens Schindel]



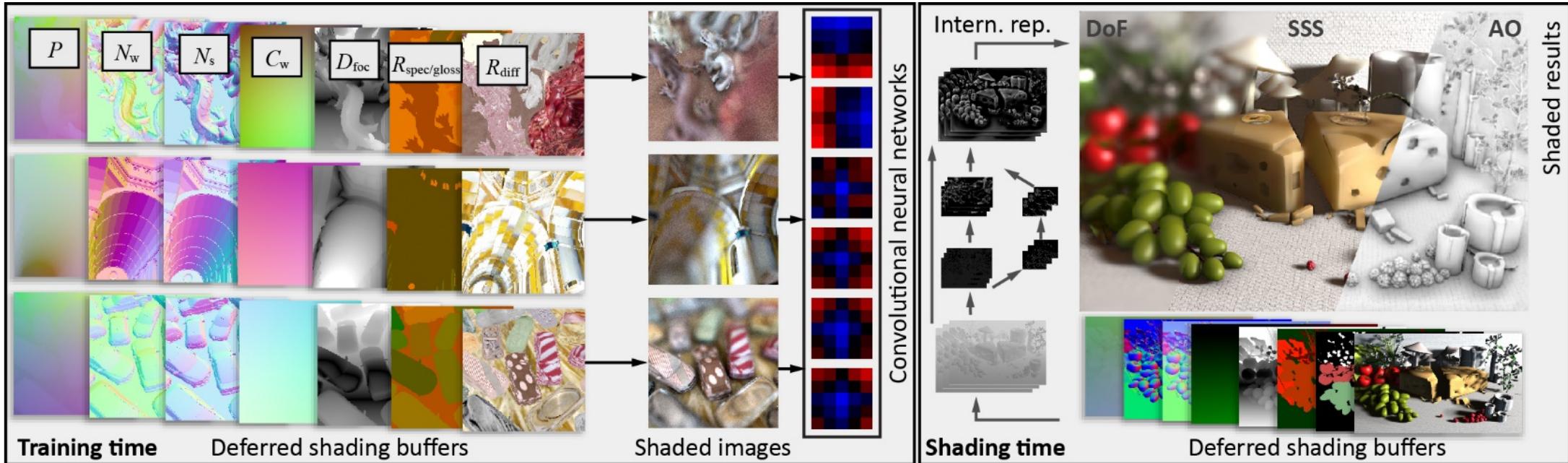
# Video Results





# Video Results





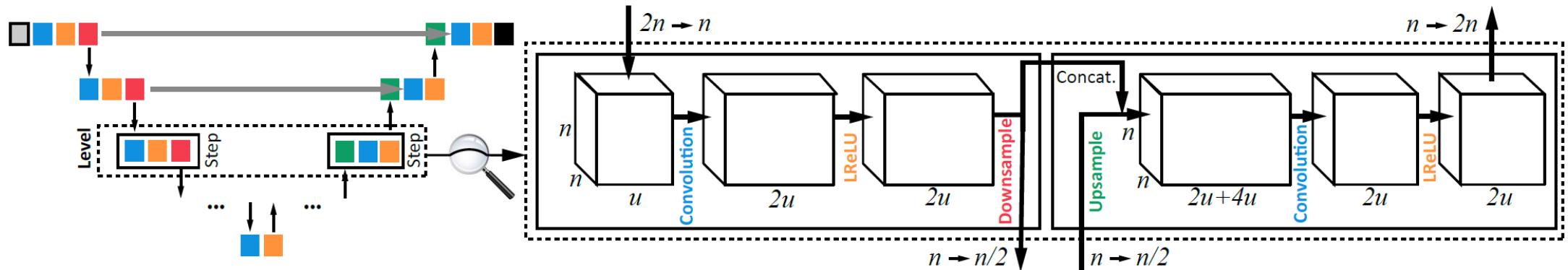
# Deep Shading – Convolutional Neural Networks for Screen-Space Shading

Nalbach et al. EGSR 2017



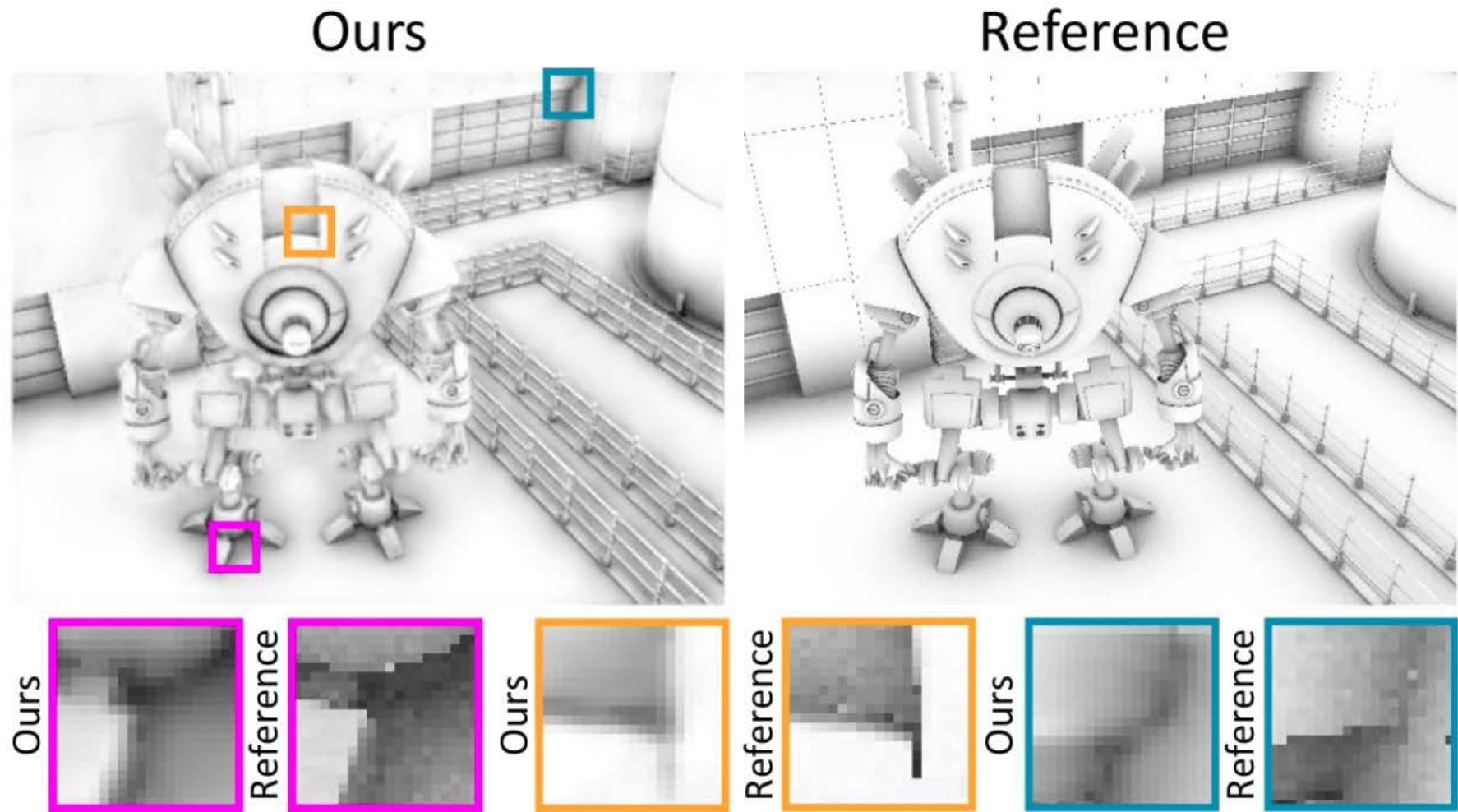
# Network Architecture

- Trained for different rendering tasks (separate training set each)
  - Ambient occlusion
  - Depth of field
  - Subsurface scattering
  - Image-based lighting
  - AO + DoF
  - Ambient occlusion
  - Depth of field
  - Subsurface scattering
  - Image-based lighting
  - AO + DoF
  - ...



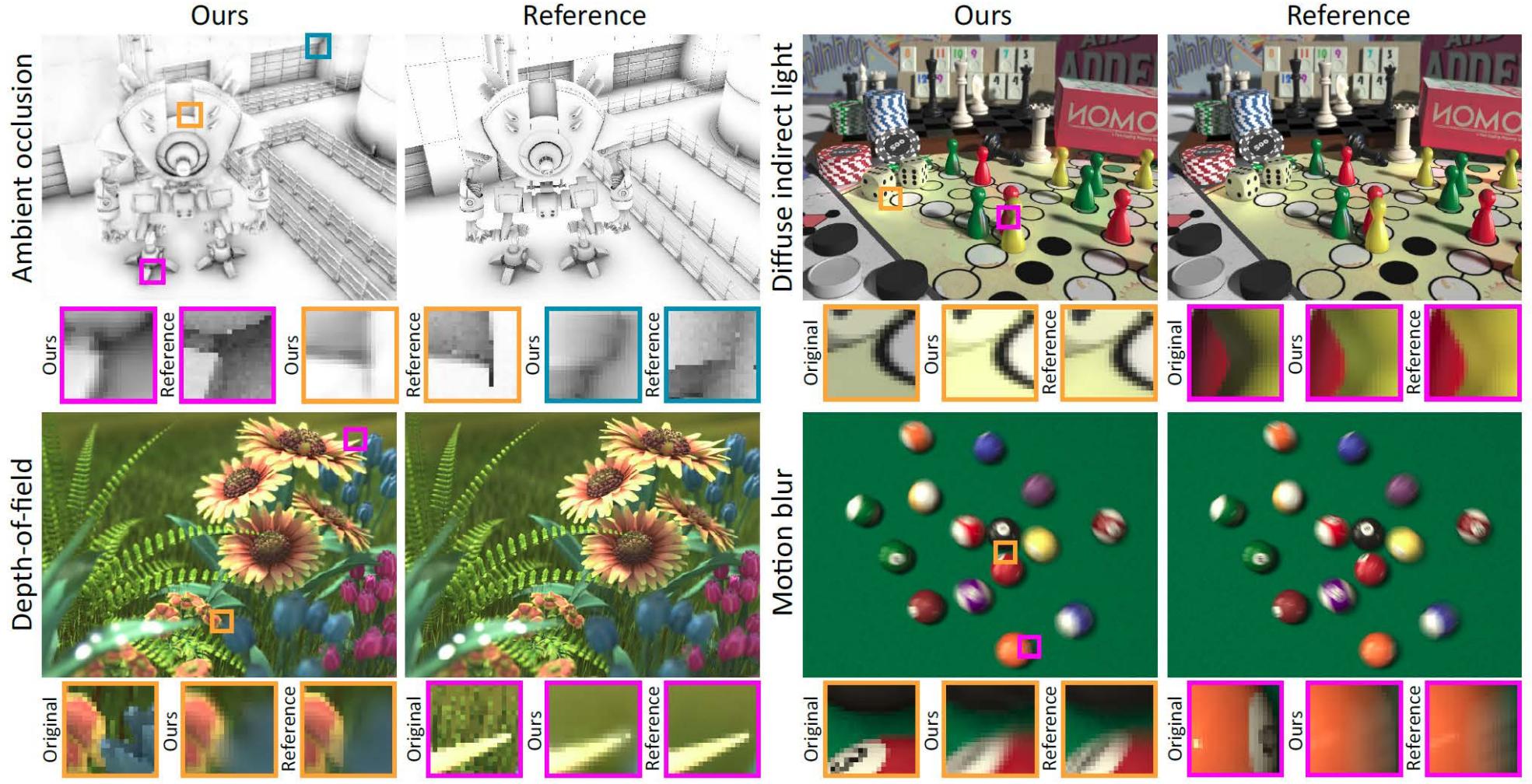


# Results





# Results



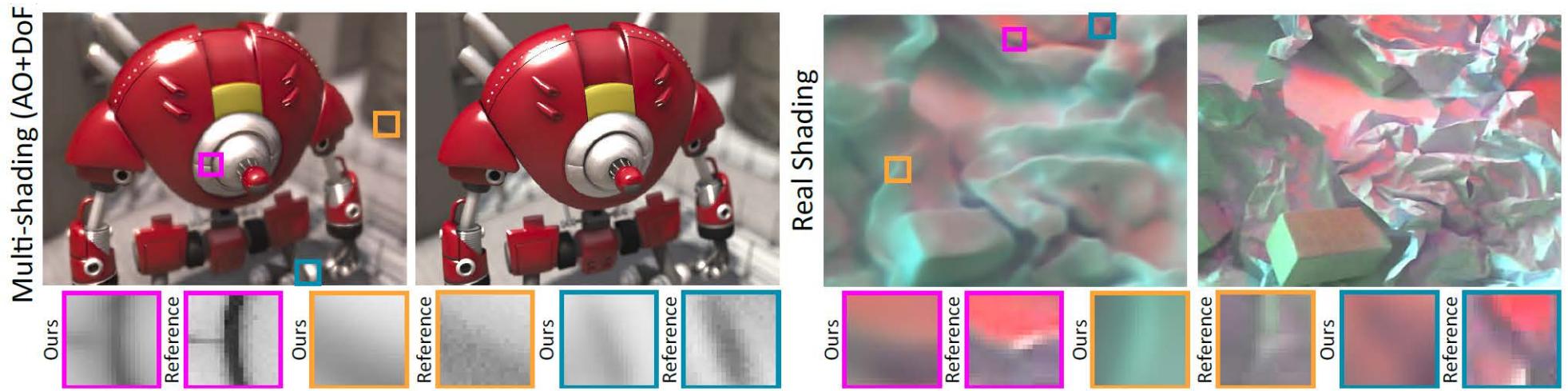


# Results





# Results



## References

---

- Patrick Wieschollek, Ido Freeman, Hendrik P.A. Lensch, Learning Robust Video Synchronization without Annotations, ICMLA 2017
- Nima Khademi Kalantari, Steve Bakó, Pradeep Sen, A Machine Learning Approach for Filtering Monte Carlo Noise, ACM SIGGRAPH 2015
- Vogels et al., Denoising with Kernel Prediction and Asymmetric Loss Functions, ACM SIGGRAPH 2018.
- Daniel Holden, Jun Saitoy, Taku Komuraz, Neural Network Ambient Occlusion, SIGGRAPH Asia 2016
- Oliver Nalbach and Elena Arabadzhiyska and Dushyant Mehta, Hans-Peter Seidel and Tobias Ritschel, Deep Shading – Convolutional Neural Networks for Screen-Space Shading, EGSR 2017