

Inference of scene representations for video analysis

Georgios Georgiadis

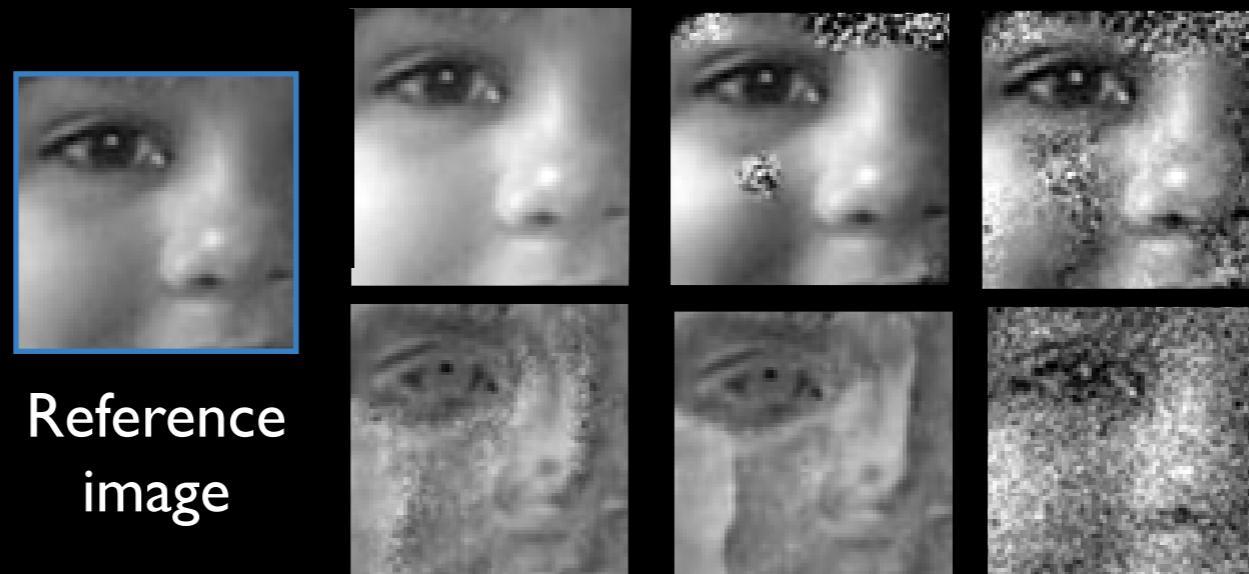
Video compression: a quick introduction

- **GOAL:** Reduce redundancy in video data
- **Main techniques:** Spatial image compression and motion compensation
- **Lossless compression schemes** is not necessary for video compression
- **Lossy compression** can produce **perceptually indistinguishable images**
- Trading off perceptual similarity with encoding bits

Perceptual similarity

- PSNR is ineffective in measuring perceptual similarity

All noisy images have the same PSNR value



Wang et al., "Mean Squared Error: Love it or Leave it?", IEEE Signal Processing Magazine, 2009

- Assumption: We are given a metric that correctly measures the perceptual dissimilarity between two images or two videos
- What is the highest compression rate we can achieve for a given fidelity (in terms of perceptual similarity)

Modeling of pixel values

- Current encoders model pixel values

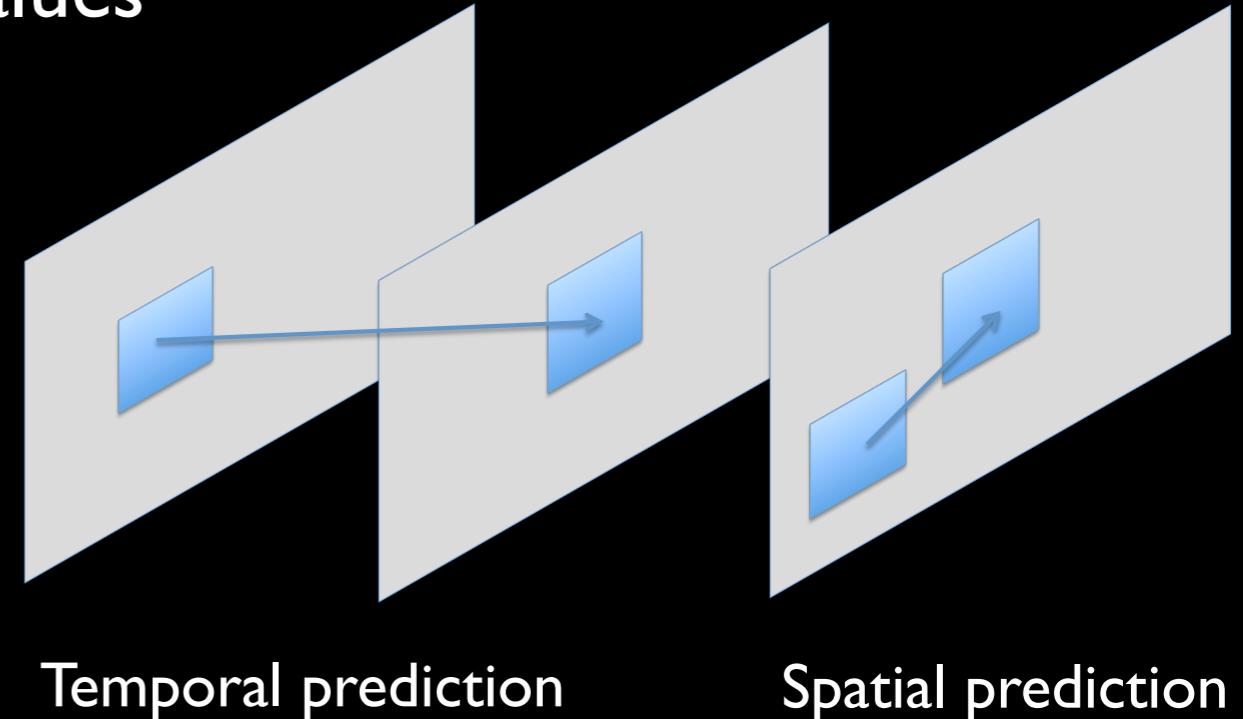
- Successful because:

1. high compression ratio
2. fast decoding

- Problem: Hit compression wall

- CISCO projects:

1. By 2016: 1.2 million video minutes (>2 years) will travel the Internet every second
2. Fastest consumer mobile service will be mobile video (from 271 million in 2011 to 1.6 billion in 2016)
3. By 2016 video traffic will be 86% of global consumer traffic



Modeling the scene

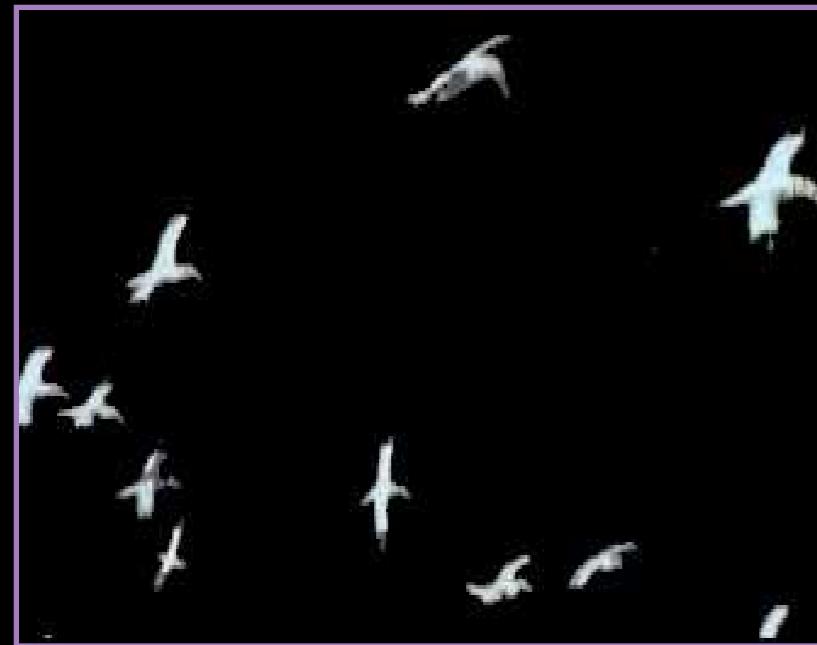
- Key idea: Modeling the pixel values is not necessary to achieve perceptual indistinguishability
- The scene is an object that has a certain:
 - Photometry (appearance)
 - Dynamics (motion)
 - Geometry (shape)
 - Topology (occlusions)
 - Statistical regularities (e.g. textures)
- Encoder should pay attention to visual phenomena the user is most sensitive (e.g. occlusions)
- Modeling the scene gives greater flexibility on how to spend encoding bits (compared to modeling the pixel values)
- Consequence: Ignore details that the viewer cannot perceive

Main approach

Texture / structure partition



Visual Structures



Temporal redundancy



Visual Textures



Spatial redundancy

Visual Textures



Regular textures



Stochastic
textures



Domain-deformed
textures



Range-deformed
textures

- Definition: A **texture** is a region of an image that can be rectified into a sample of a stochastic process that is:
 - **Locally stationary** (defining characteristic)
 - **Ergodic** (allows inference from a single realization)
 - **Markovian** (allows finite dimensional model)

Visual textures

Stationarity and ergodicity

Notation:

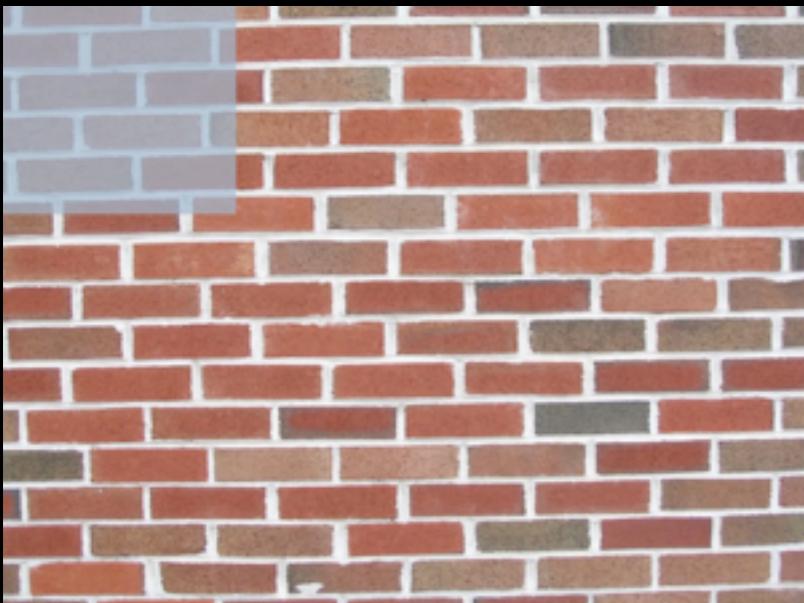
$$I \sim dP(I)$$

$$\phi : \mathbb{R}^N \rightarrow \mathbb{R}^W \text{ Statistic}$$

$$\phi_\omega \doteq \phi(I(\omega)) \text{ "Local" Statistic}$$

$$\omega \subset \mathbb{Z}^2$$

$$I(\omega) \doteq \{I(x), x \in \omega\}$$



Stationary process

$$\mathbb{E}(\phi_\omega) = \mathbb{E}(\phi_{\omega+T}), \quad T \in \mathbb{R}^2$$

Stationary and ergodic process

$$\frac{1}{N} \sum^N \phi_{\omega+T_i} \xrightarrow{\text{a. s.}} \mathbb{E}(\phi_\omega)$$

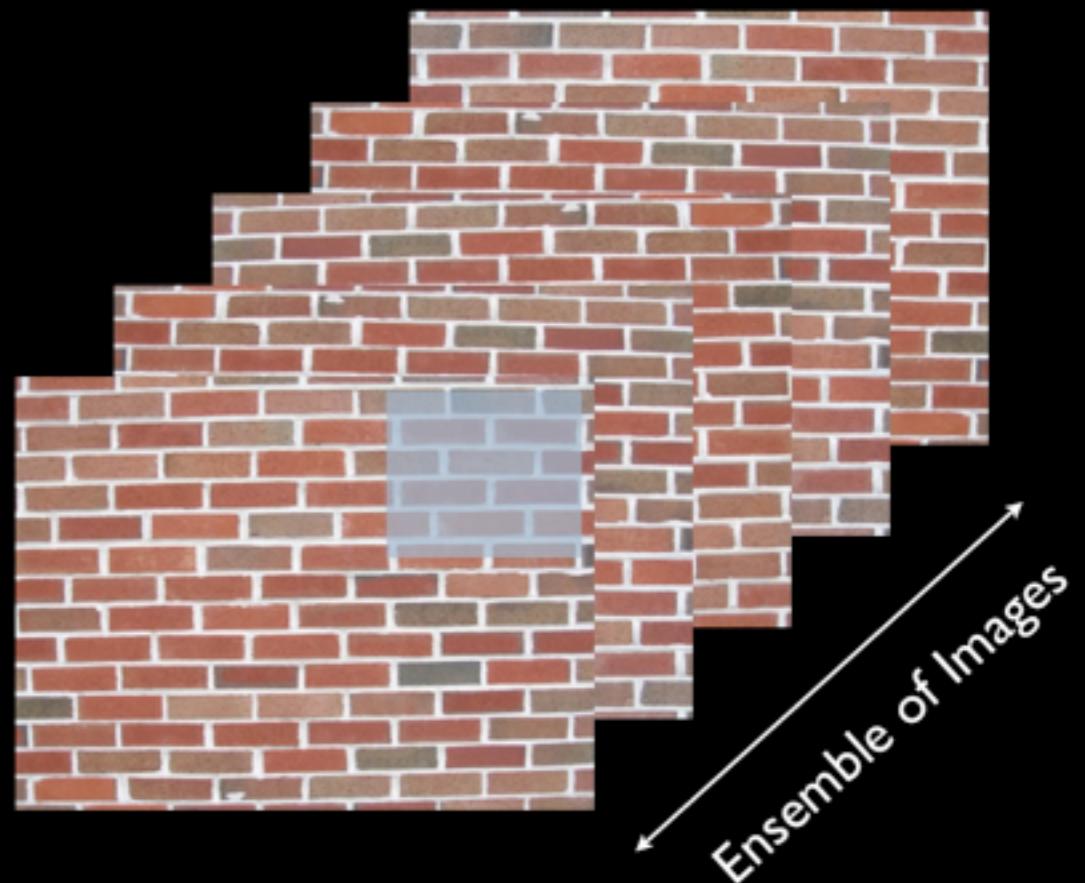
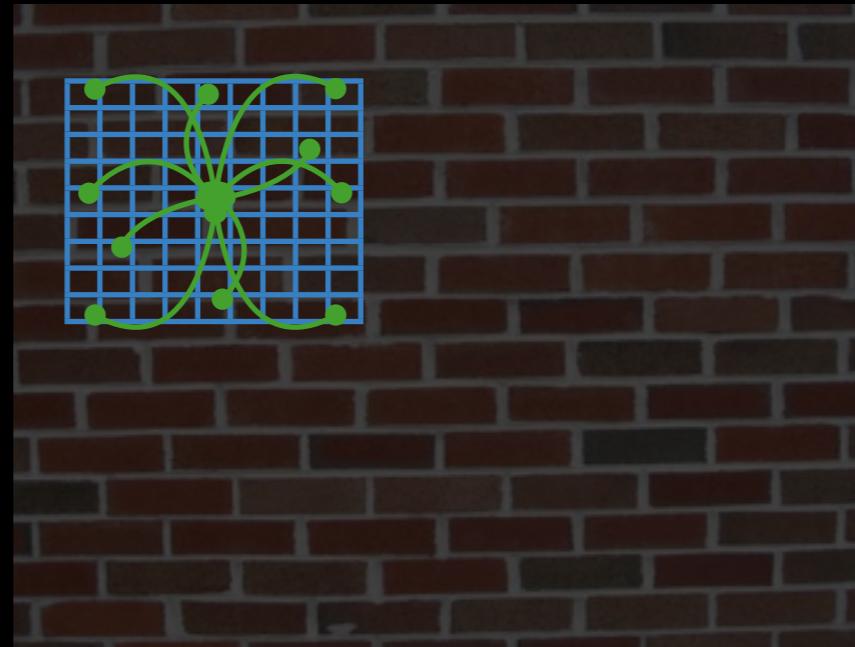
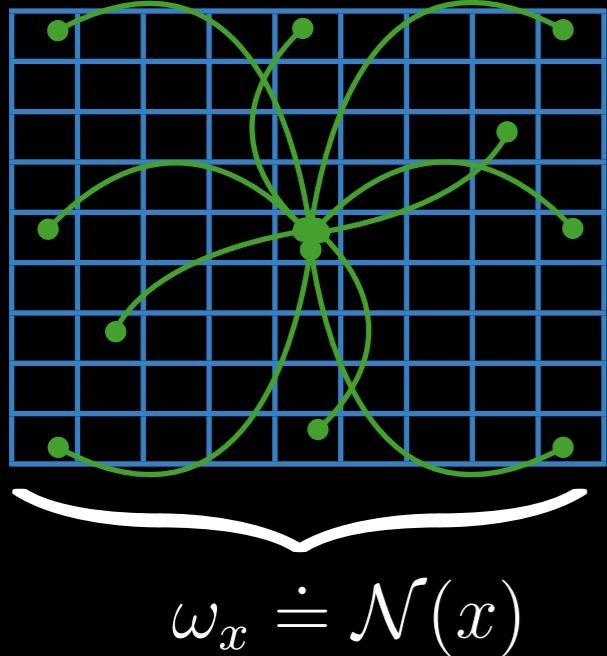


Figure caption:

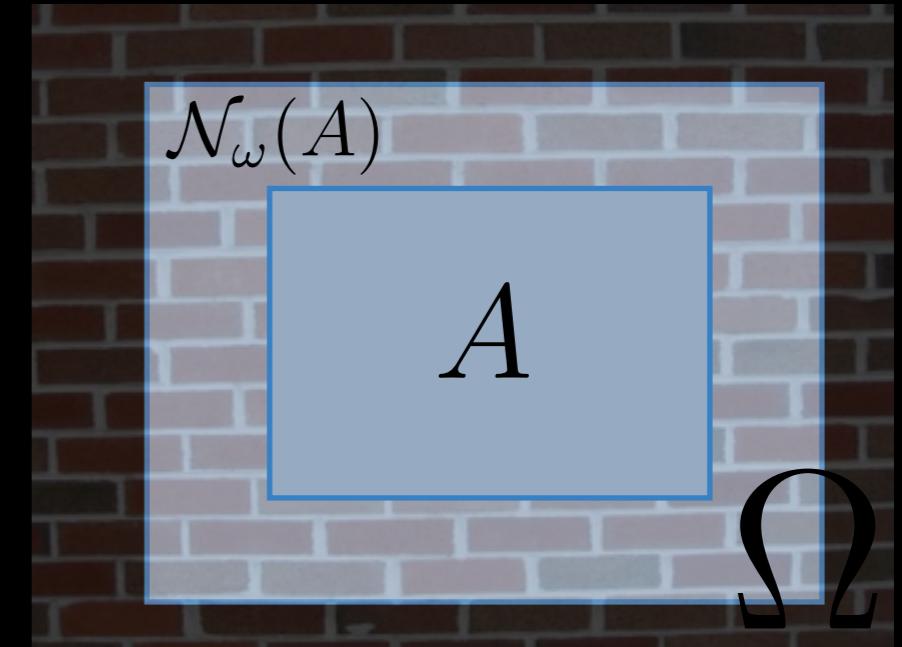
$$\phi_\omega = I(\omega) \text{ (Grayed moving square)}$$

Visual Textures: Markovianity

Neighborhood
structure of x



Homogeneity of neighborhood
structure of x



A and its neighborhood $N(A)$

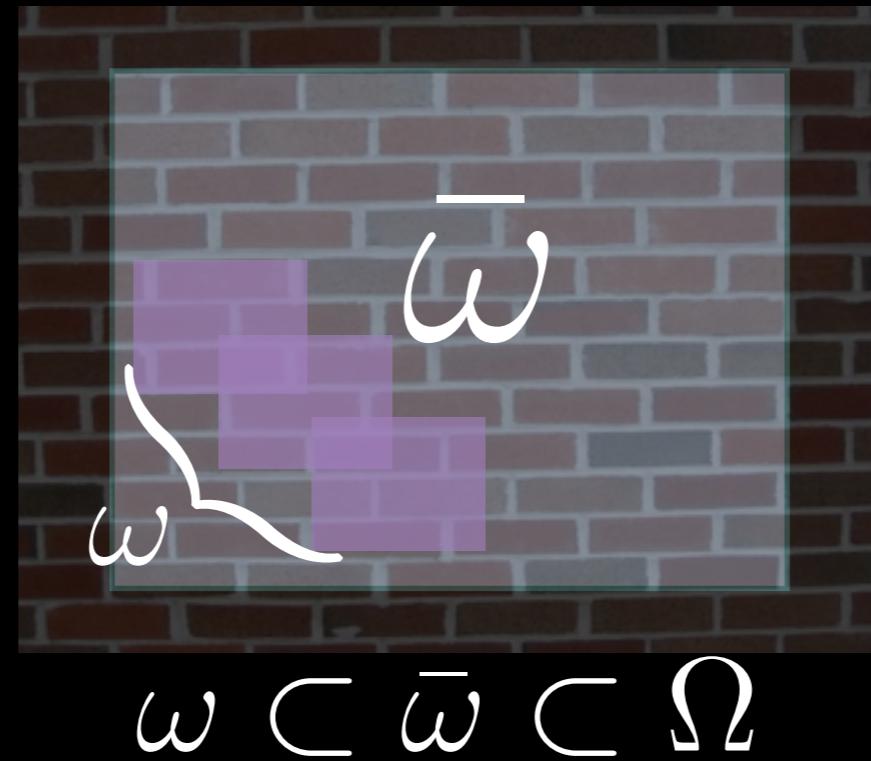
$$I(A) \perp \underbrace{I(\Omega - A)}_{\text{"outside"}}, \underbrace{I(N_\omega(A))}_{\text{"inside"}}, \forall A \subset \Omega$$

$$H[I(A)|I(N_\omega(A))] = H[I(A)|I(\Omega - A)]$$

$$\hat{\omega}(\beta) = \arg \min_{\omega} H[I(A)|I(N_\omega(A))] + \frac{1}{\beta} |\omega|$$

Visual Textures Representation

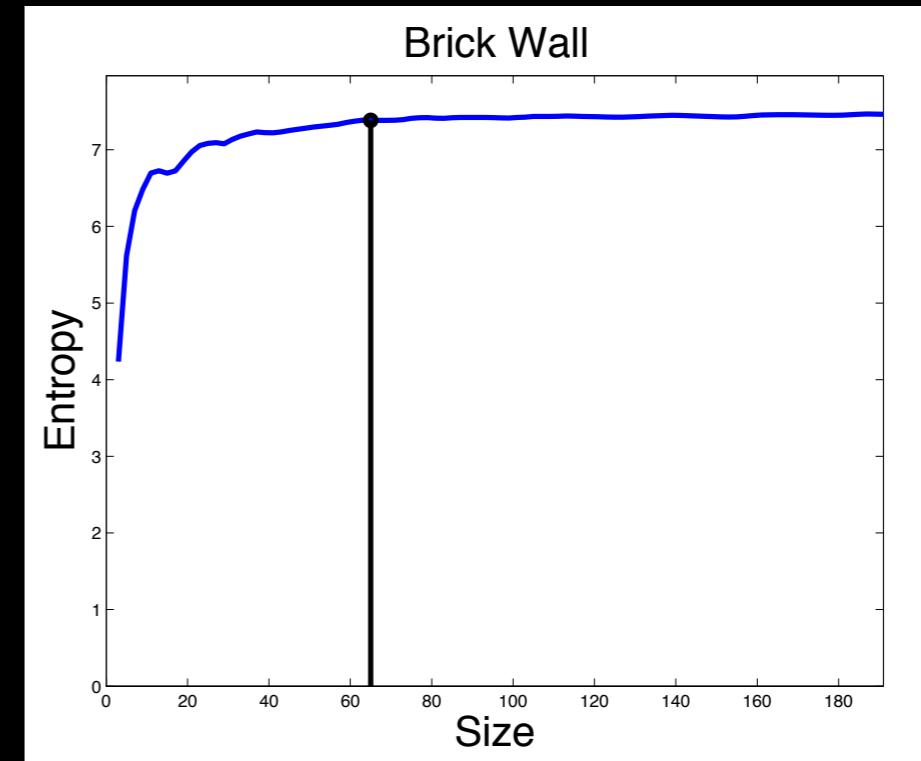
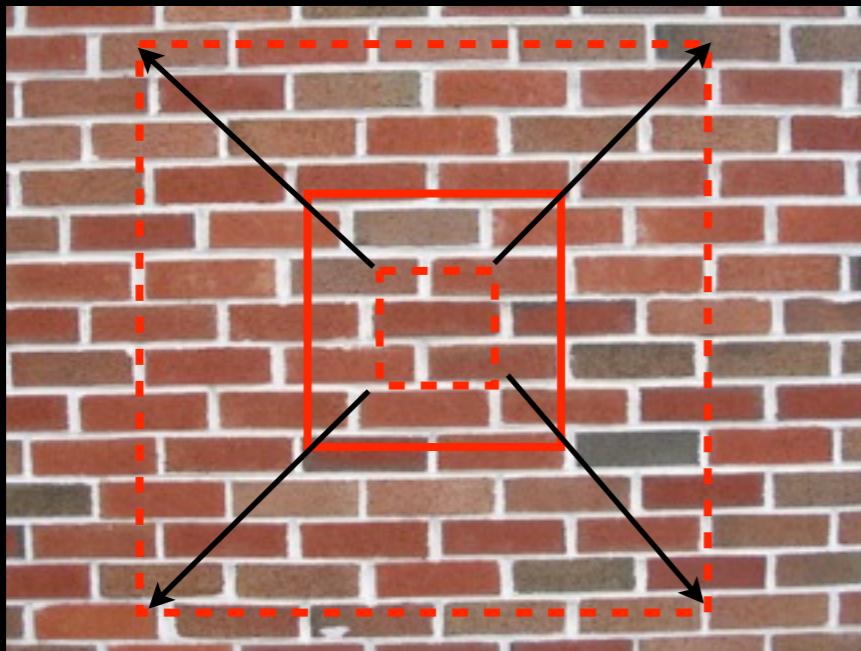
$$H [I(A)|I(\mathcal{N}_\omega(A))] \simeq -\frac{1}{M} \sum_{i=1}^M \log p [I(A_i)|I(\mathcal{N}_\omega(A_i))]$$
$$A_i \doteq A + T_i \subset \Omega$$



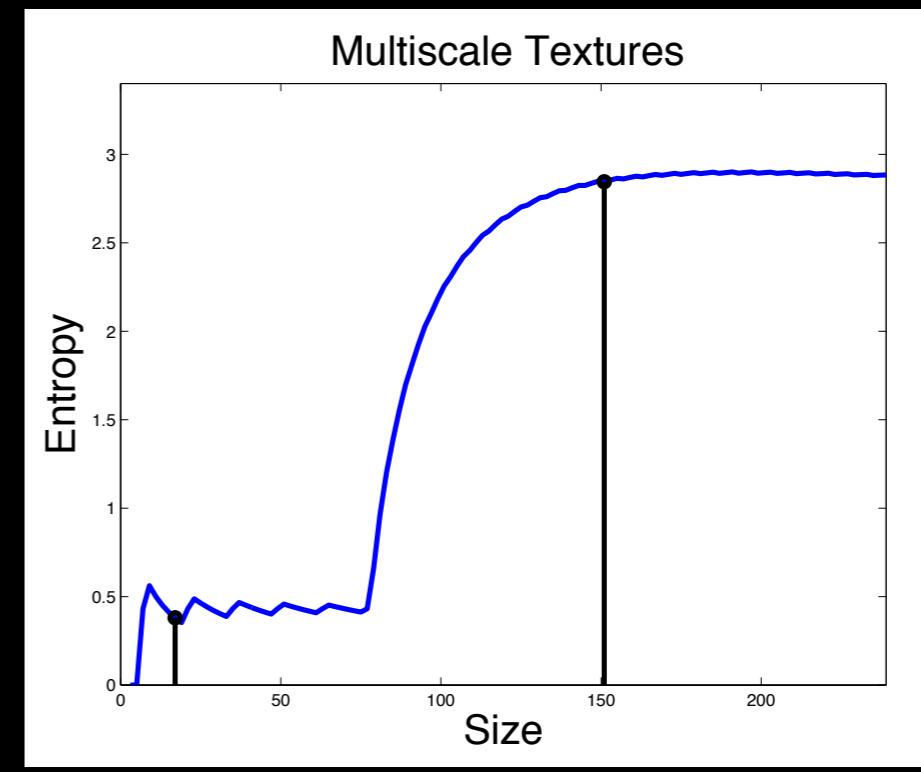
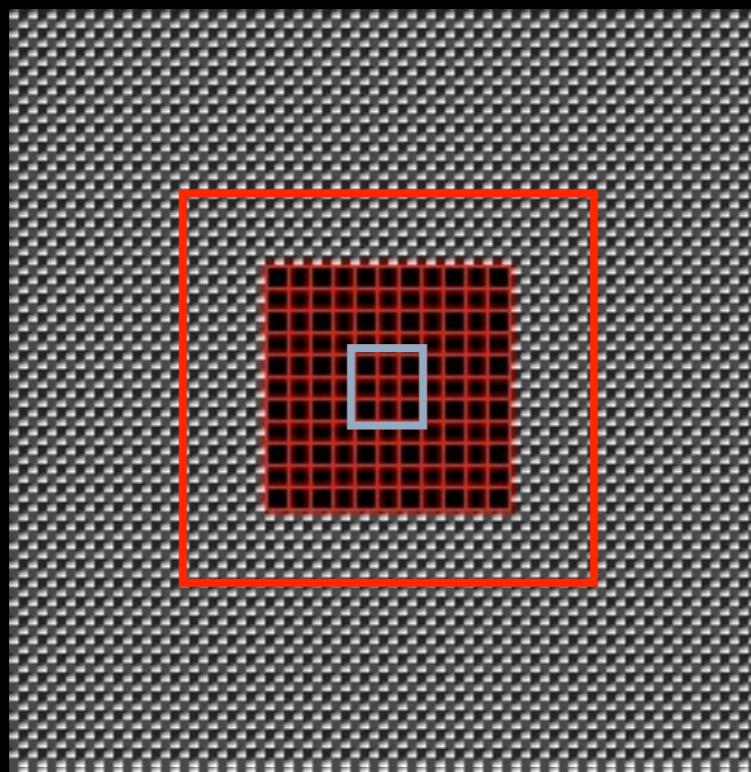
$$\hat{r}, \hat{\sigma} = \arg \min_{r=|\omega|, \sigma=|\bar{\omega}|} \hat{H} [I(A)|I(\mathcal{N}_\omega(A))] + \frac{1}{\beta} |\bar{\omega}|$$

Visual Textures

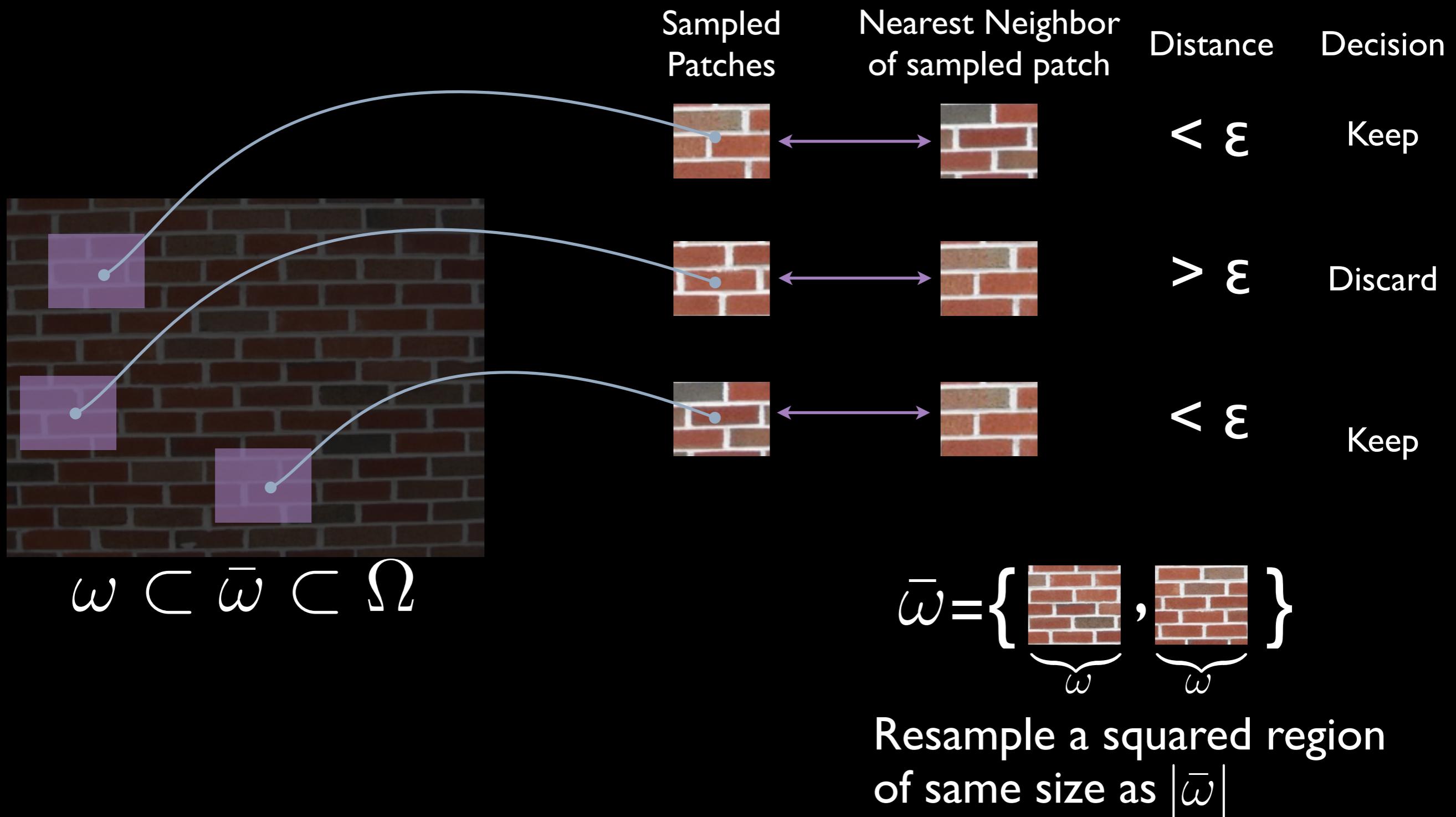
Inferring ω



Entropy of histogram of pixel values



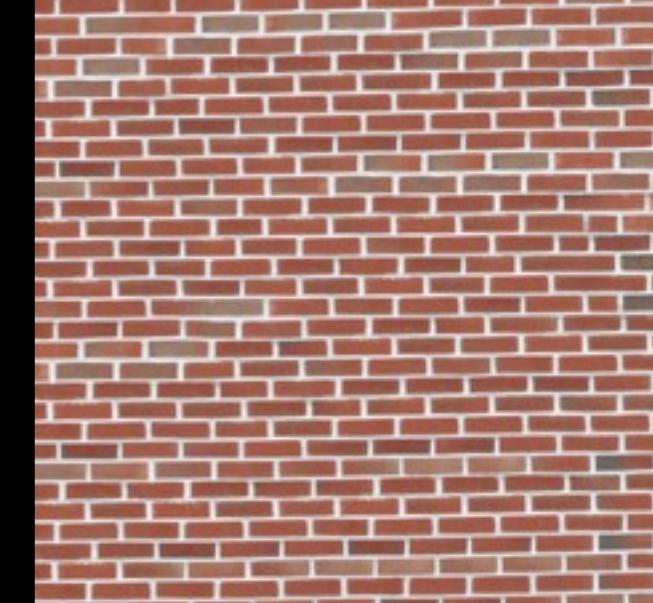
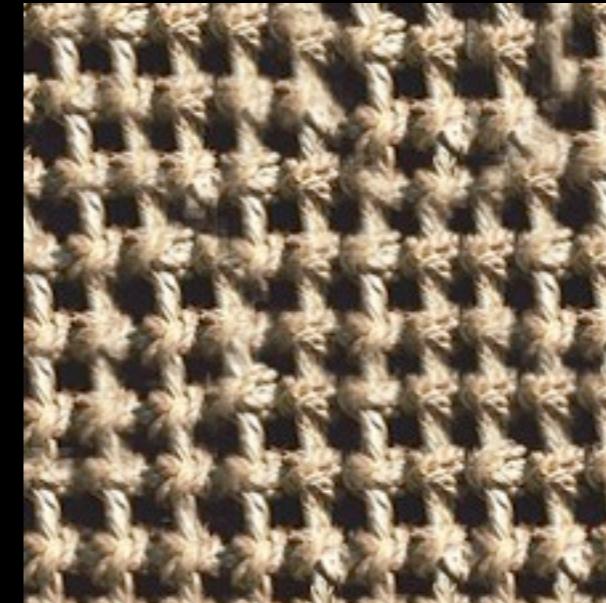
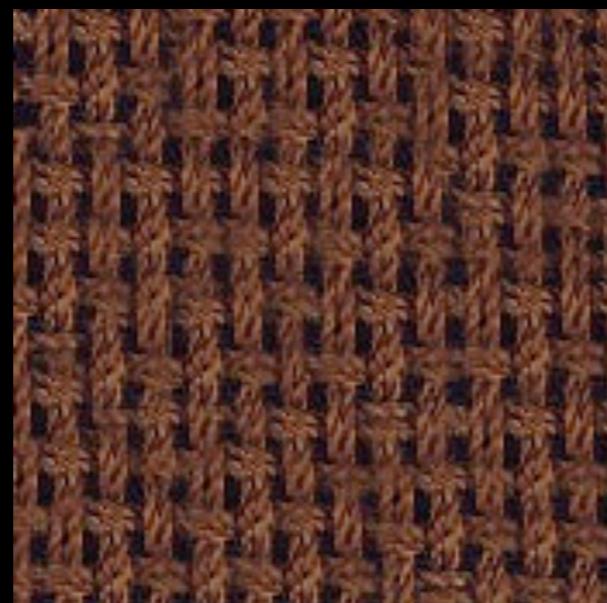
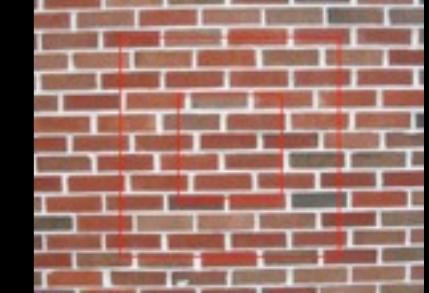
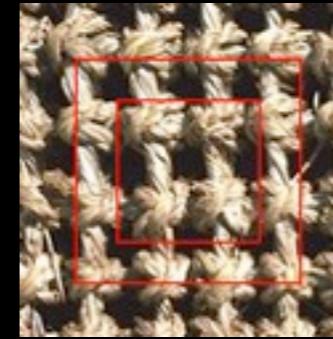
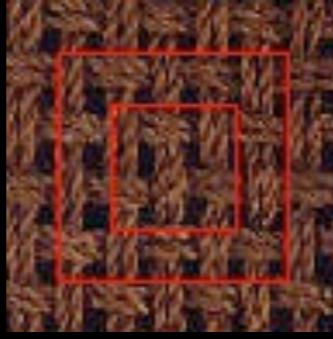
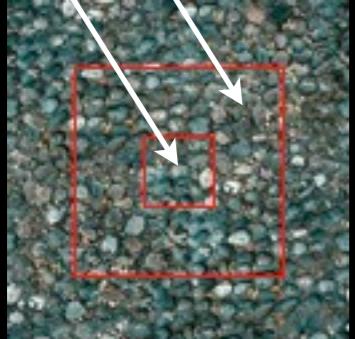
Visual Textures Inferring $\bar{\omega}$



Visual Textures

Inferring representations - Results

$$\omega \bar{\omega}$$



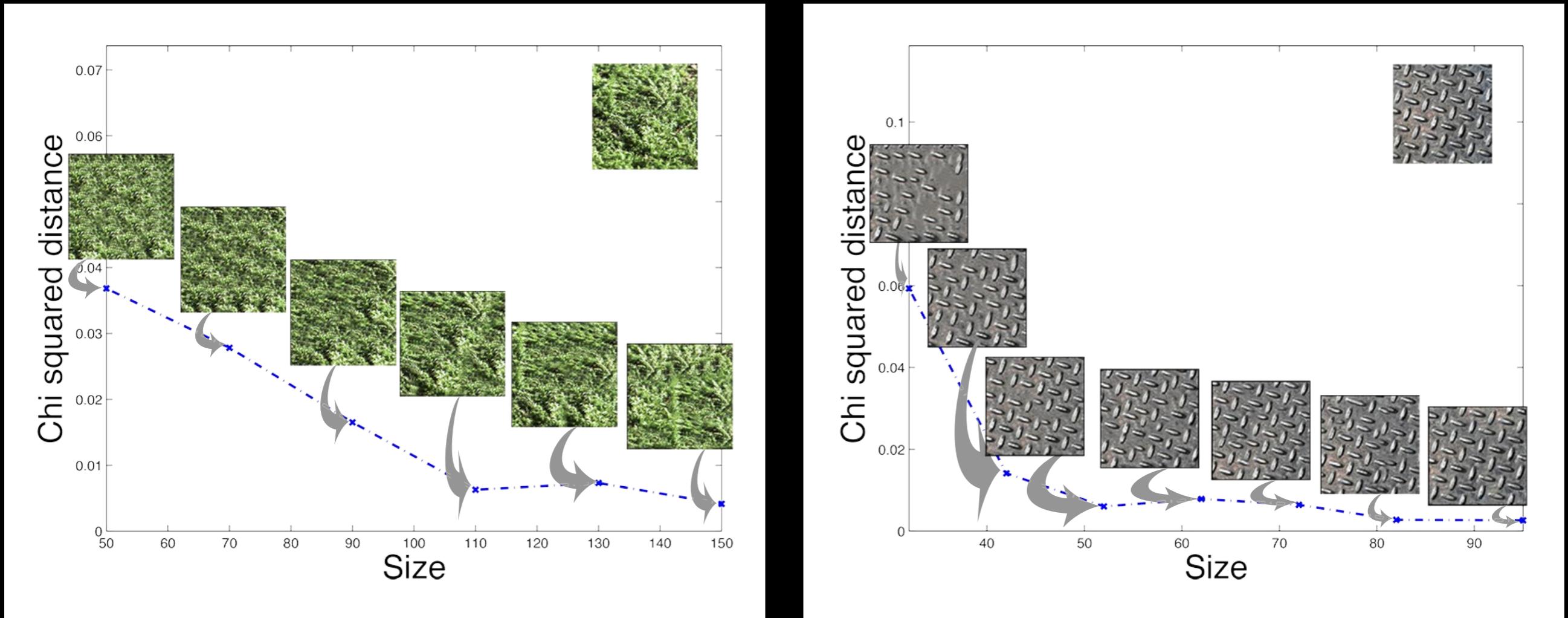
Synthesized textures from the representation

$$E(\hat{I}_s; I(\bar{\omega})) = \sum_{\tilde{\omega} \in \hat{\Omega}} \|\hat{I}_{\tilde{\omega}} - I_{\tilde{\omega}}\|^2$$

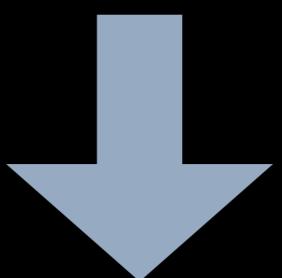
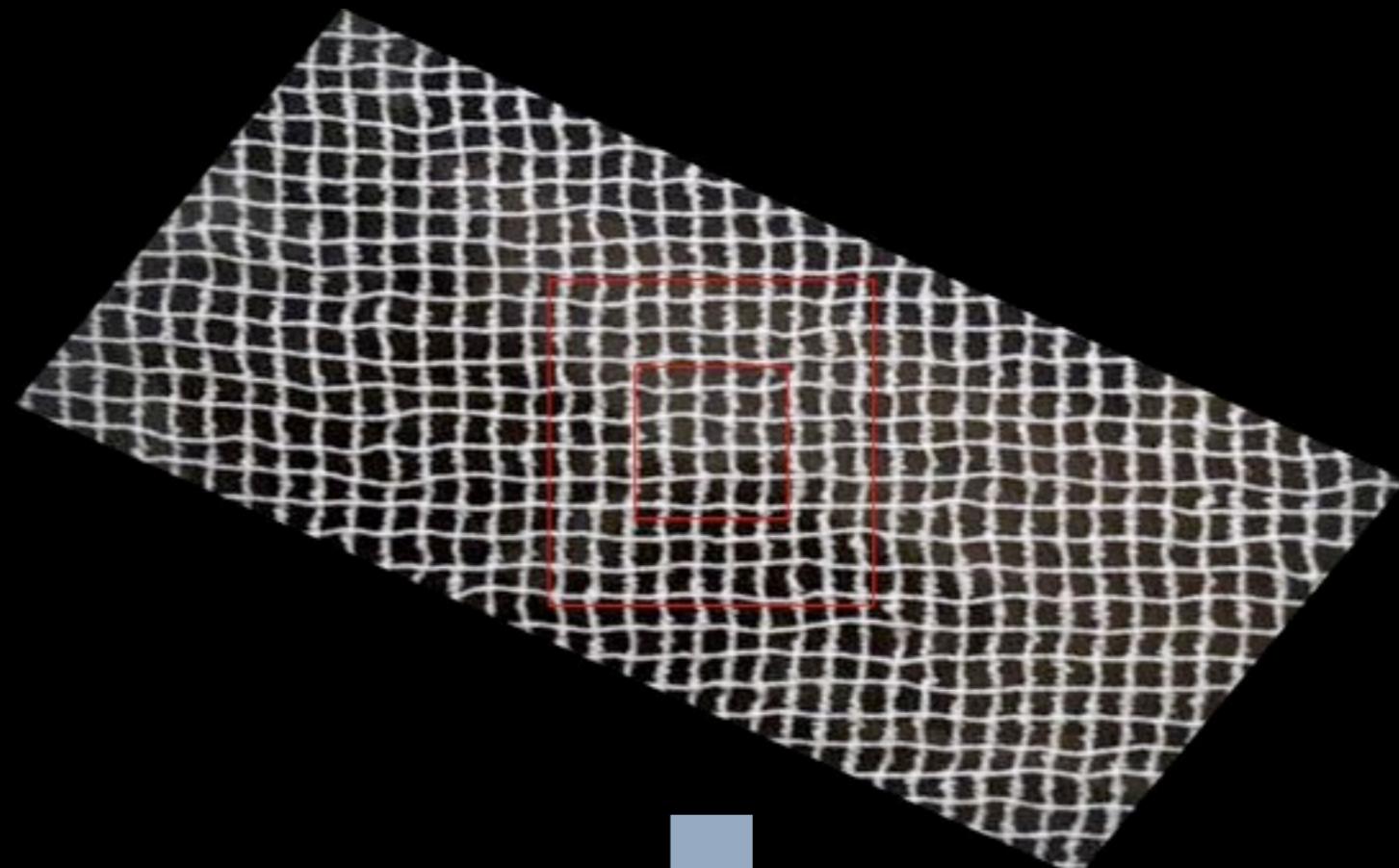
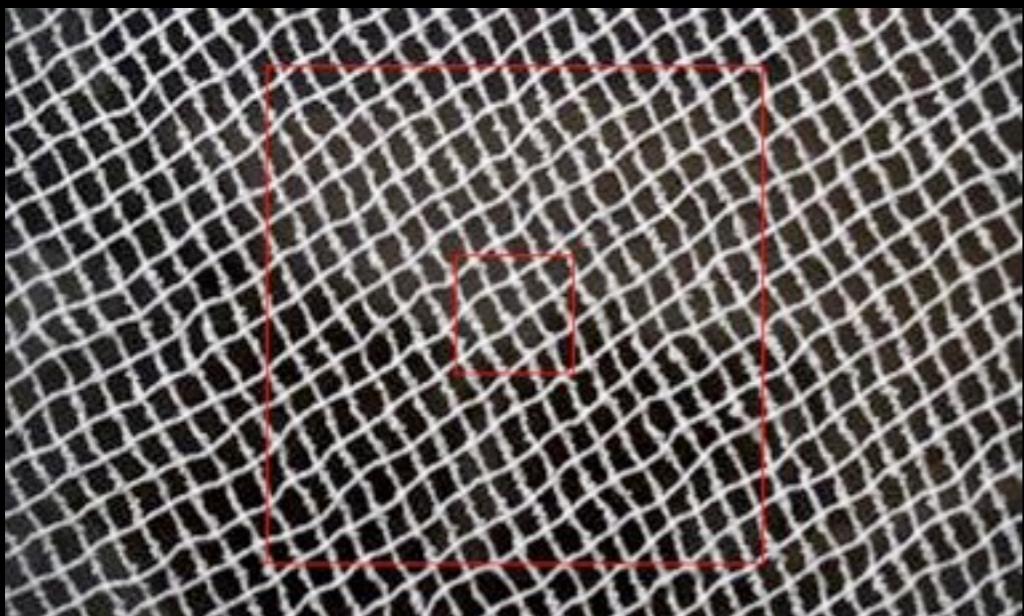
Synthesized texture Texture representation Local neighborhood $\tilde{\omega}$ in synthesized texture Nearest neighbor of $\tilde{\omega}$ in the representation

Visual Textures

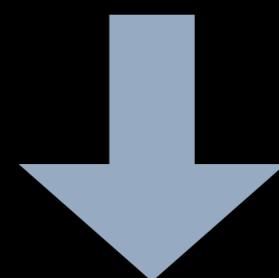
Rate - Distortion plots



Visual Textures Transformations



$$|\bar{\omega}_{rectified}| \simeq \left| \frac{\bar{\omega}_{original}}{4} \right|$$



Synthesized texture



Synthesized texture

Visual Textures Segmentation



Algorithm:

- Image is divided into small regions
- Sequentially merge regions that are the most “similar” (their pixel value histograms are the most similar)

Visual Structures

Visual Structures

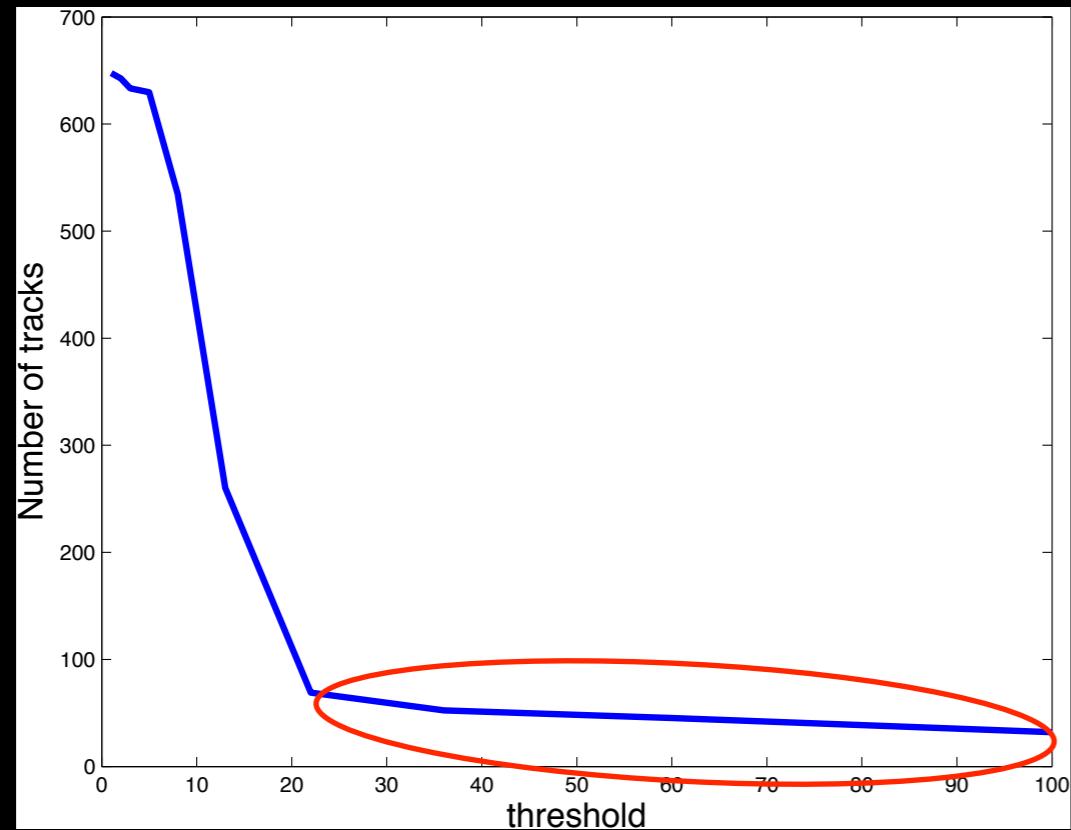


Temporal redundancy

- Structures are regions in images that:
 - trigger isolated responses of a feature detector
 - persist in more than one image
- Structures \longleftrightarrow Saliency

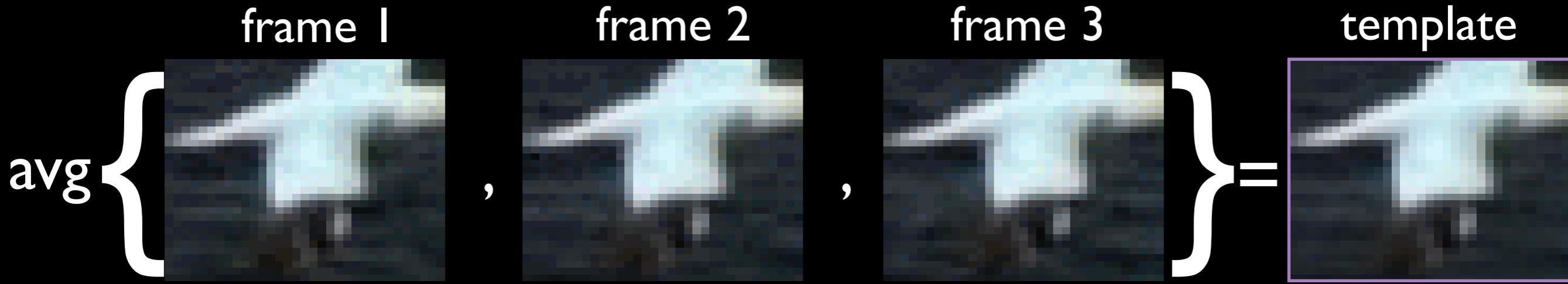
Visual Structures

Detection, tracking and selection

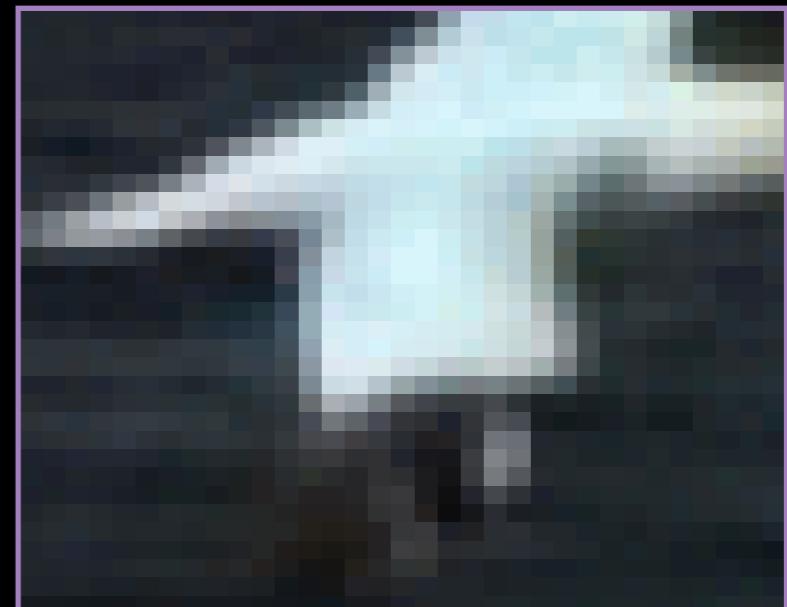
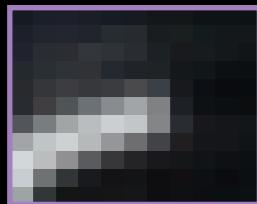


- Most reliable tracks persist over a wide range of threshold values:
 - ▶ stable reference frame (high quality)
 - ▶ large temporal support (highly compressible)

Visual Structures Representation



Template examples:

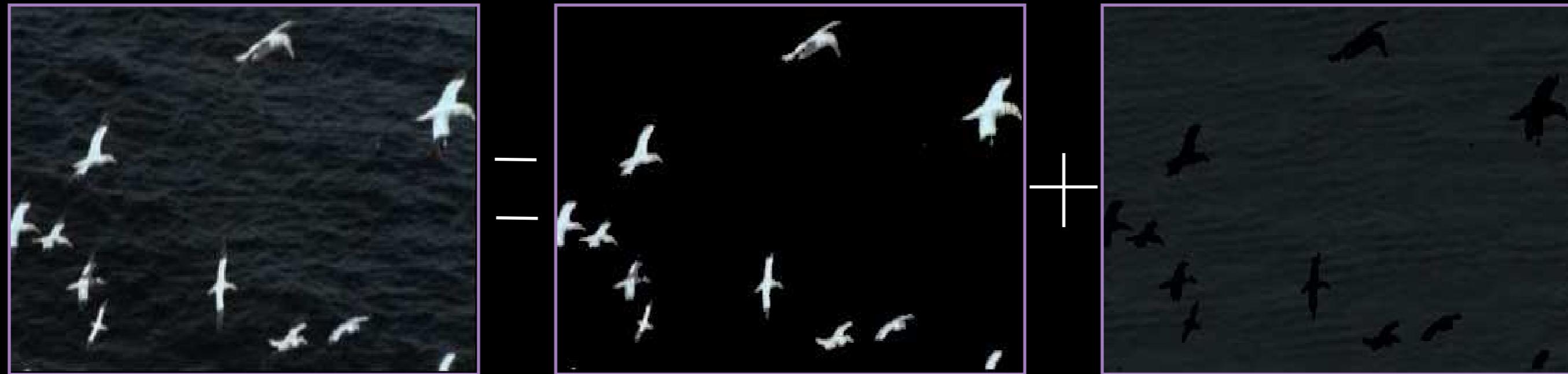


- Dictionary is composed of templates that have different spatial resolutions and are functions of different temporal supports
- Final representation: template and pixel location of central pixel

Texture / Structure partition

Visual Structures

Visual Textures



- Texture / Structure partition depends on multiple frames
- Texture / Structure partition depends on scale

Finer Scale: Structure

- For any given scale σ of observation is either a **structure** or a **texture**.

Isolated Extremum at scale ϵ

=> Not Stationary

=> Not Texture

Coarser Scale: Texture



with scale $\bar{\sigma}$ = σ

Stationary at scale $\sigma \gg \epsilon$

=> No isolated extremum at scale σ

=> Not a structure

Texture / Structure partition

A first attempt



Input frame



Structures



Texture segmentation



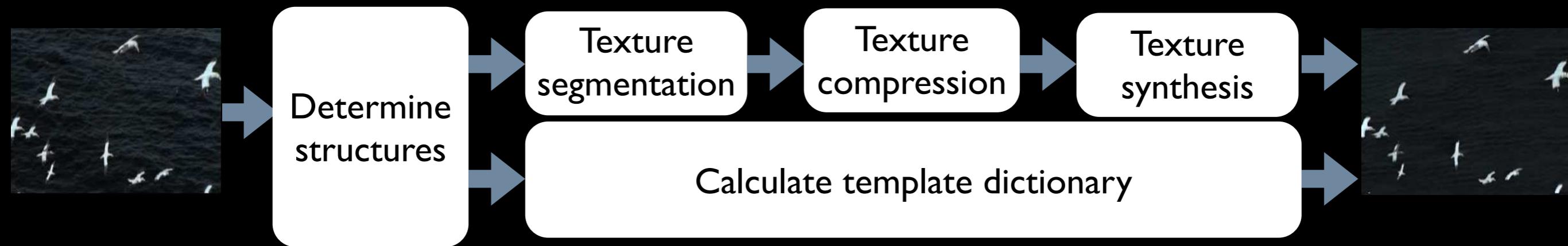
Texture compression

Texture synthesis



Reconstructed frame

Modeling the scene



Region type	Redundancy	Calculate	Representation
Structure			Template, position
Texture			Compressed representation, texture boundary

Texture / Structure partition

Qualitative results



Original video



Reconstructed video

Comparison with other methods

- Our method successfully captures salient regions of the video



Original



Zhi et al. [5] -
ICCV'11

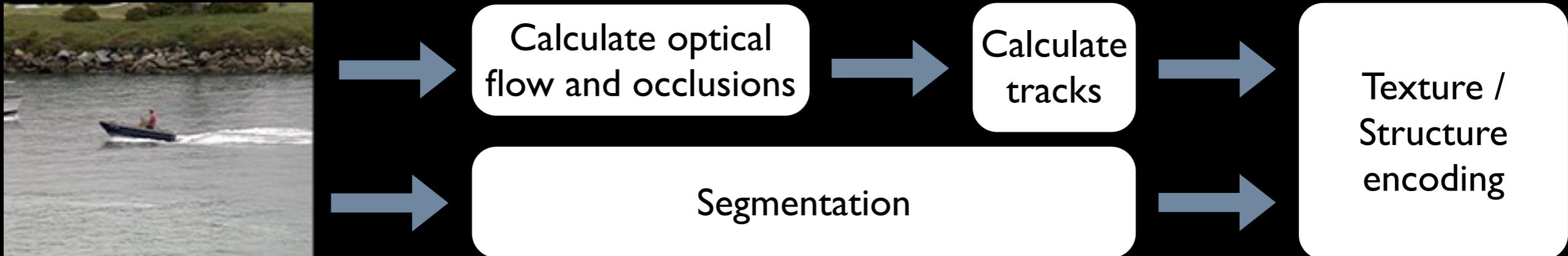


Our method

Computational cost

- For a 332x192 resolution video (7 frames):
 - Tracking (once for whole video): ~2s
 - Texture segmentation (per frame): ~1m
 - Texture compression (once per texture): ~45s
 - Texture synthesis (once per texture): ~3m
- Overall:
 - Encoding time per frame: ~66s
 - Decoding time per frame: ~30s (texture synthesis)

Modeling the scene (a year before)



Region type	Redundancy	Calculate	Encode
Structure	 → 	Velocities	Median motion vector
Texture		Crop / Downsample	
Occlusions/ Homogeneous	?		Intensities

Temporal / Spatial / Textural partition



Input frame



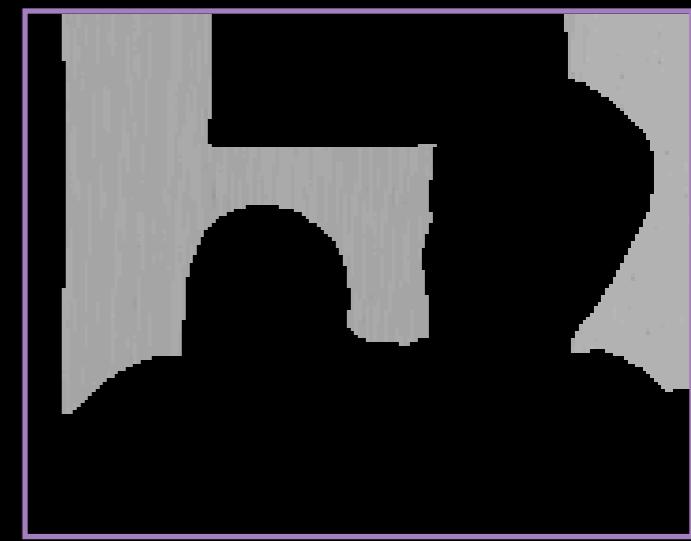
Temporal prediction

+



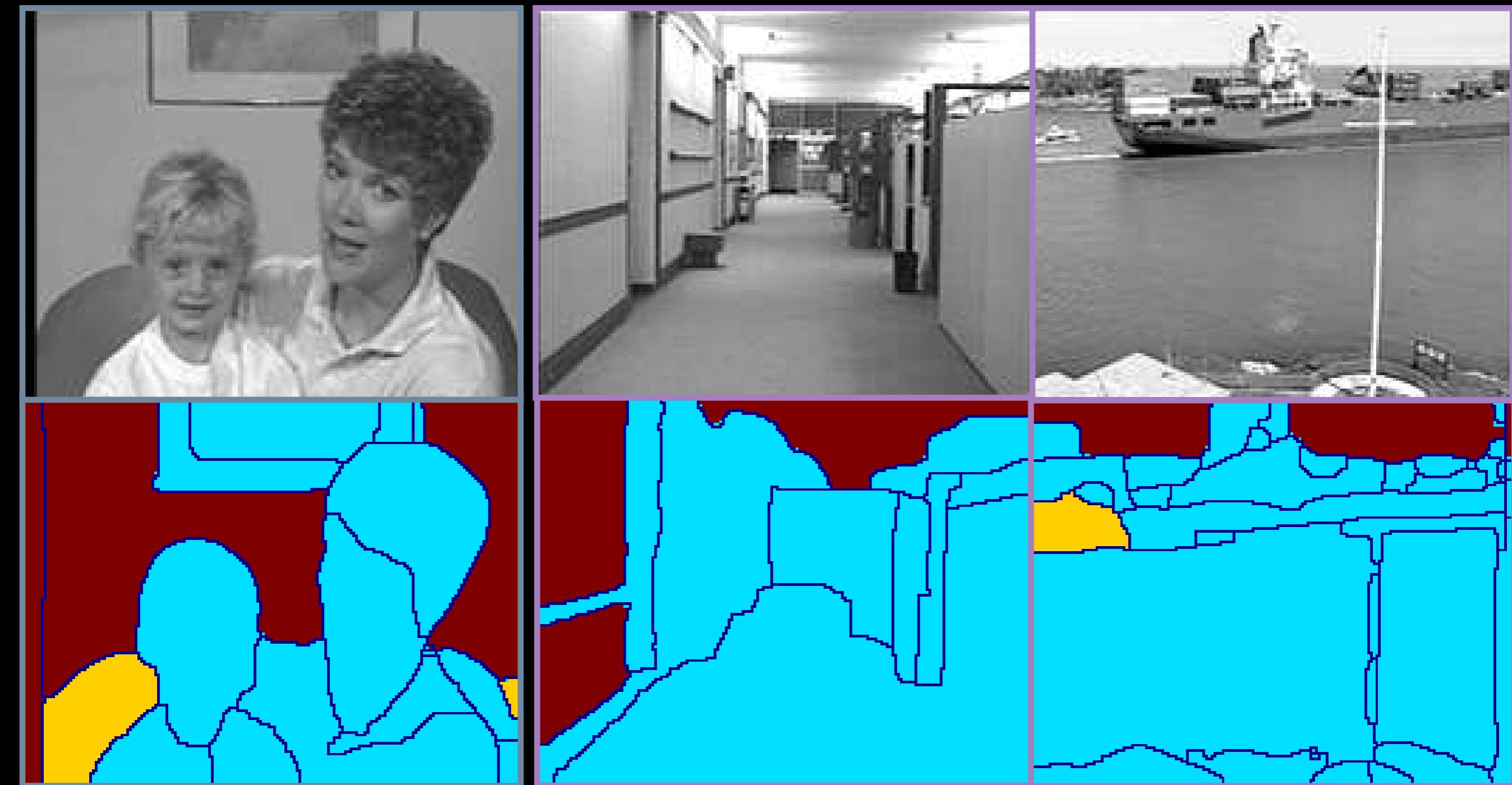
Spatial prediction

+



Texture synthesis

Partition results



Blue: Temporal prediction

Yellow: Spatial prediction

Red: Texture synthesis

Failure cases

Qualitative results



Predicted frame

Error frame

Reconstructed frame



Original video



Our method

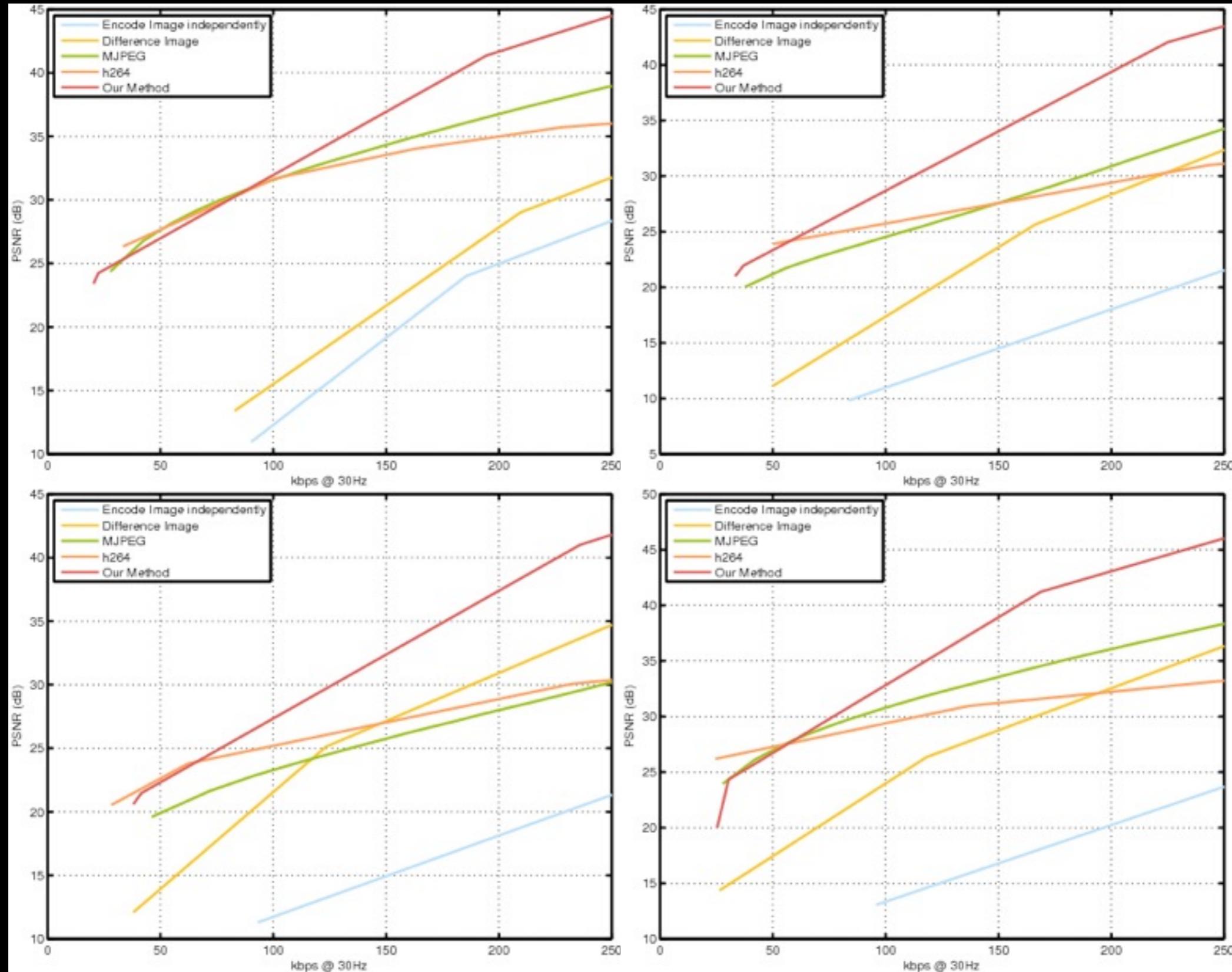
Same
file
size

Reconstructed videos



H.264

Quantitative results



Computational cost

- For a VGA resolution video:
 - GPU implementation of optical flow: ~10s
 - Tracking: ~2s
 - GPU implementation of segmentation: ~13s
 - Texture synthesis: ~5s
- Segmentation can run in parallel with optical flow and tracking
- Encoding time per frame: ~21s

Actionable Saliency Detection



Occlusion detection:

- Detect salient regions as violations of co-visibility hypothesis:
$$H0 = \{\exists \text{ a diffeo } w : \Omega \rightarrow D \mid I(x, t + dt) - I(w(x), t) \stackrel{\text{IID}}{\sim} \mathcal{N}\}$$
- No violation when moving over homogeneous backgrounds

Approach

Co-visibility hypothesis:

$$H0 = \{\exists \text{ a diffeo } w : \Omega \rightarrow D \mid I(x, t + dt) - I(w(x), t) \stackrel{\text{IID}}{\sim} \mathcal{N}\}$$

Rigidity:

$$H1 = \{\exists T \in \mathbb{R}^3, \omega \in \mathbb{R}^3, Z : D \rightarrow \mathbb{R}^+ \mid w(x) = \pi(\exp(\hat{\omega})\bar{x}Z(x) + T)\}$$

Notation:

$$x \in \mathbb{R}^2 \quad \text{Pixel location}$$

$$y(x) \in \mathbb{R}^2 \quad \text{Motion vector}$$

$$V \in \mathbb{R}^3 \quad \begin{matrix} \text{Translational} \\ \text{velocity of} \\ \text{camera} \end{matrix}$$

$$\omega \in \mathbb{R}^3 \quad \begin{matrix} \text{Rotational} \\ \text{velocity of} \\ \text{camera} \end{matrix}$$

$$Z : D \rightarrow \mathbb{R}^+ \quad \text{Depth map}$$

$$y(x) = \mathcal{A}(x) \frac{V}{Z(x)} + \mathcal{B}(x)\omega$$

$$\underset{V}{\text{minimize}} \quad \psi(V) = \frac{1}{2} \|\hat{C}(V)Y\|_2^2.$$

$$[I - \mathcal{W}C(C^T\mathcal{W}^T\mathcal{W}C)^{-1}C^T\mathcal{W}^T] \mathcal{W}Y \doteq \hat{C}(V)Y = 0.$$



Results



Future work

- **Structures:** template selection

- **Textures:**

- (i) Texture segmentation (>2 regions)

- (ii) Texture synthesis: speed-up?

- (iii) Dynamic textures

- (iv) Structured textures

- **Texture / Structure partition**

Conclusions

- Model the scene rather than the pixel values
- We provided definitions, developed representations and designed algorithms for **textures** and **structures**
- Next generation video compression schemes need to model the **scene**, **semantics**, **context**. Encoding should depend on the **application**, the **viewer** and many more

Questions?