# Data Quality In The Cultural Heritage Sector: From An Image Processing Perspective

Tan Lu

Department of Mathematics and Data Science

Vrije Universiteit Brussel (VUB)

Department of Digitization

Royal Library of Belgium (KBR)

KBR

images

VRIJE UNIVERSITEIT BRUSSEL

ADOCHS

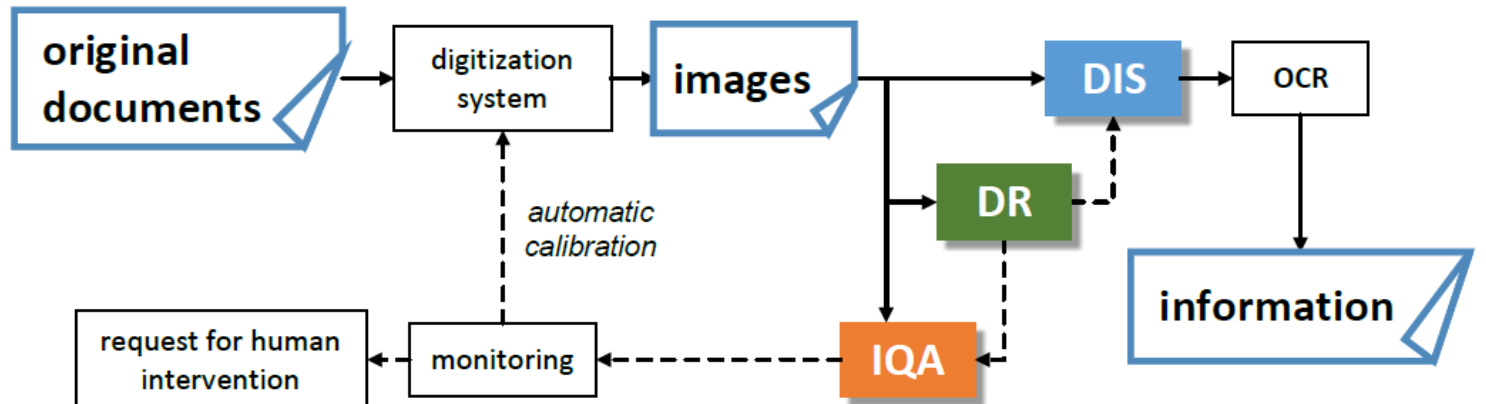.arch CEGESOMA

metadata

ULB UNIVERSITÉ LIBRE DE BRUXELLES

➢ *Digitization in the Cultural Heritage Sector*

- ❑ *Google Books*: over 40 million books
- ❑ *Europeana Newspapers*: aggregating 18 million historic newspaper pages and converting 10 million newspaper pages to full text
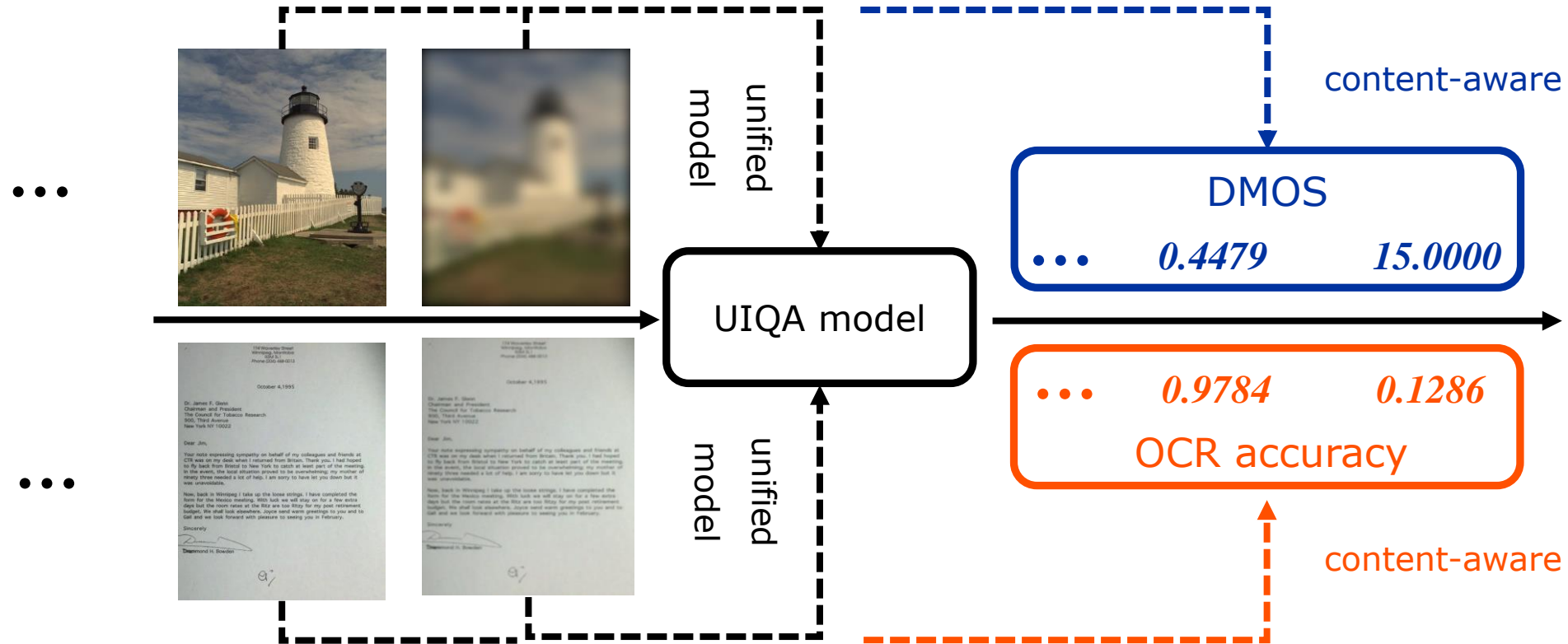- ❑ *Royal Library of Belgium (KBR)*: 4500 medieval codices and about one million prints and drawings

➢ *Challenges and Opportunities*

- ❑ *Image Quality Assessment (IQA)*
- ❑ *Image Understanding*
  - • Document Image Segmentation (DIS)
  - • Damage Recognition (DR)

# A Unified Approach to Image Quality Assessment



- A unified model to process natural and document images simultaneously

- Content-aware such that different types of quality information is provided according to different types of input images

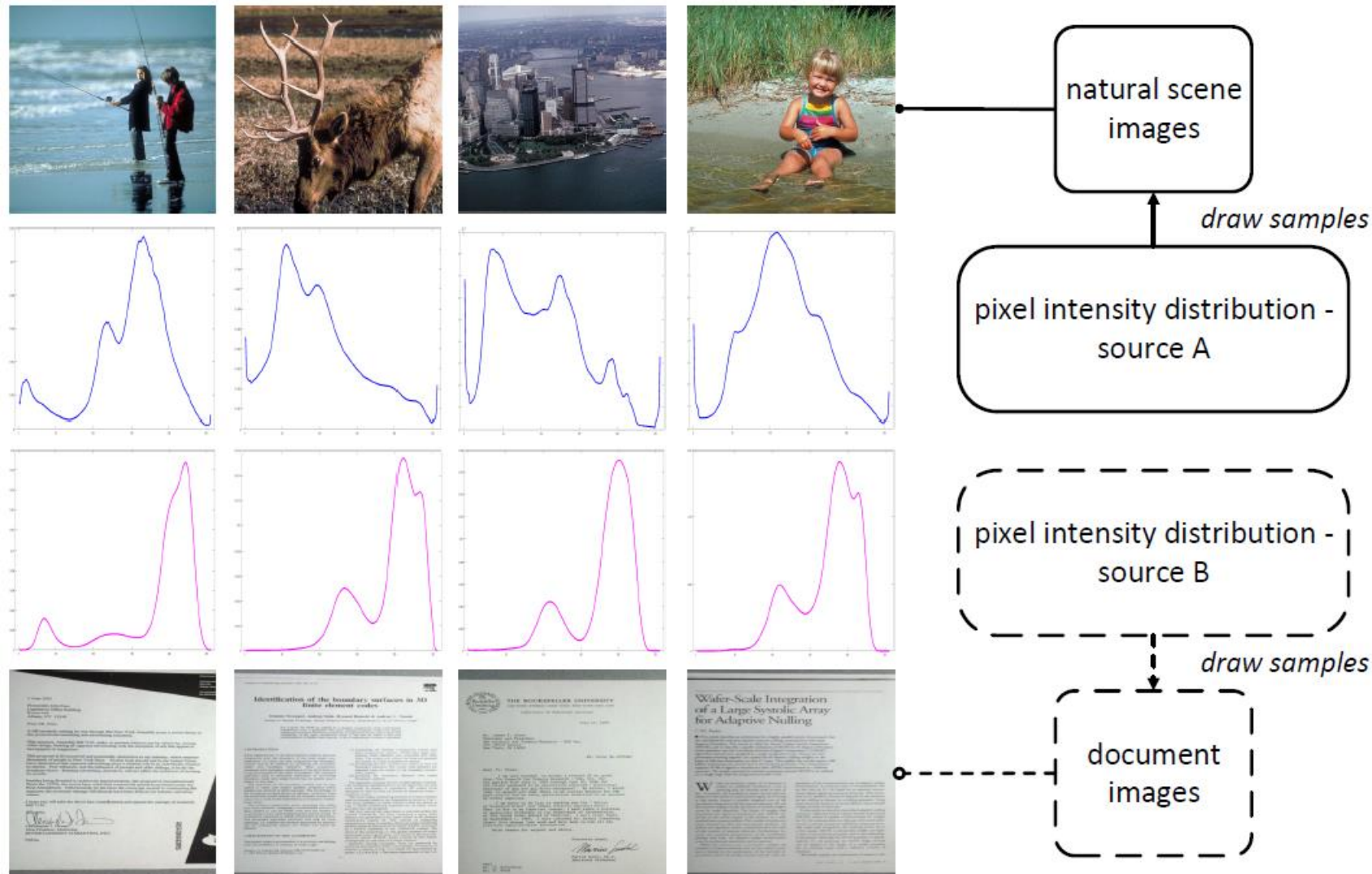# Document Image Quality Assessment based on Transfer Learning

# Document Image Quality Assessment based on Transfer Learning

| DIQA Model | PLCC | SRCC |
|---|---|---|
| CORNIA | 0.937 | 0.862 |
| CNN | 0.950 | 0.898 |
| LDA | - | 0.913 |
| HOS | 0.960 | 0.909 |
| Sparse Model | 0.935 | 0.928 |
| RNN | 0.956 | 0.916 |
| proposed method | **0.965** | **0.931** |

| DIQA Model | Document-wise | | General | |
|---|---|---|---|---|
| | PLCC | SRCC | PLCC | SRCC |
| CORNIA | 0.9747 | 0.9286 | 0.9370 | 0.8620 |
| Focus | 0.9378 | **0.9643** | 0.6467 | - |
| MetricNR | 0.9750 | 0.9107 | 0.8867 | 0.8207 |
| CG-DIQA | 0.9523 | 0.9429 | 0.9063 | 0.8565 |
| proposed method | **0.9763** | 0.9550 | **0.9651** | **0.9312** |

➢ ***Cross-Domain Homogeneity between Natural and Document Images***

The knowledge learned on natural image processing can be effectively exploited for the OCR accuracy prediction of document images.
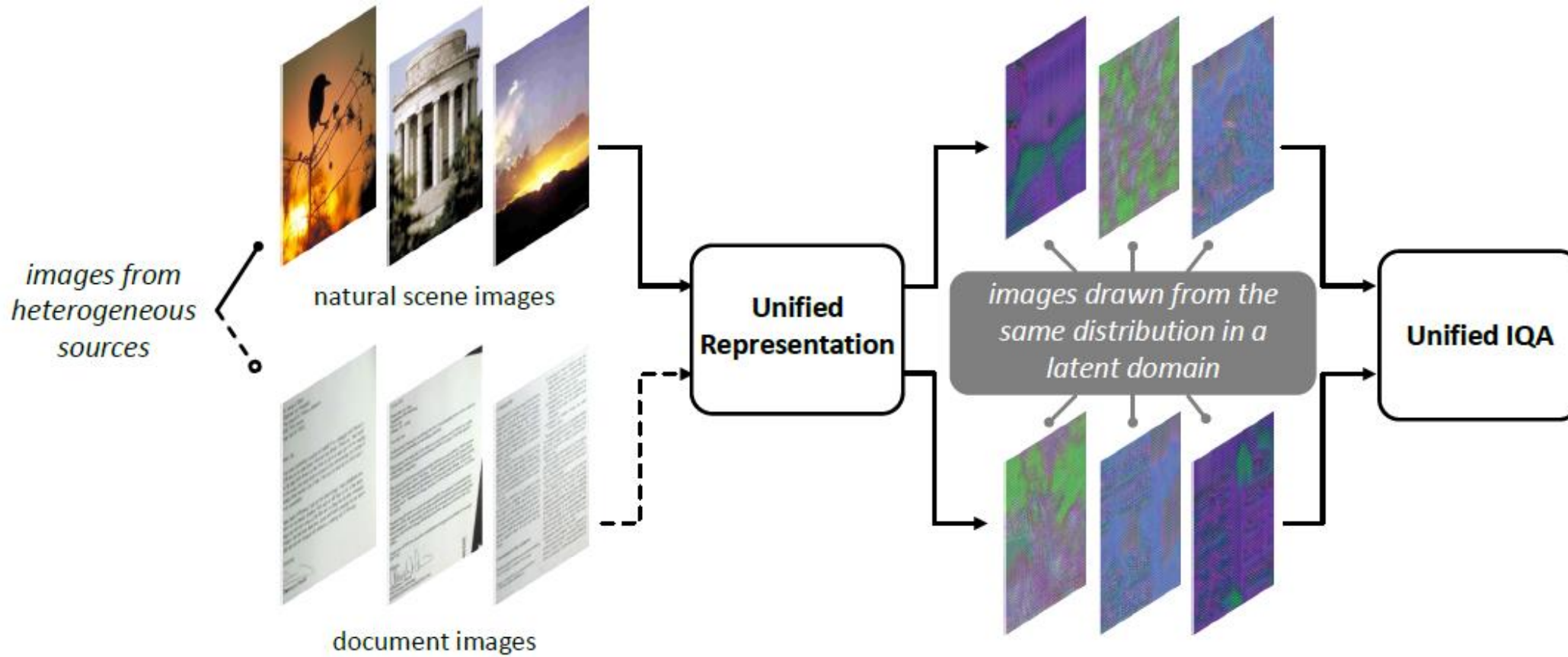
# Unified Image Quality Assessment



- natural scene images
  - *draw samples*
- pixel intensity distribution - source A
- pixel intensity distribution - source B
  - *draw samples*
- document images

➤ ***Cross-Domain Homogeneity between Natural and Document Images***

- Possible to process natural and document images simultaneously within one quality assessment model

- Balanced performance on these two types of images can be obtained with the UIQA model

- The process of learning a common representation is mixed with that of regressing the common representation towards respective quality scores – difficult to investigate and develop

# Unified Image Quality Assessment based on Contractive GAN



> *Cross-Domain Homogeneity between Natural and Document Images*

- Learning a common representation (i.e. a generalization) of natural and document images in a latent domain

- The process of generalization is separated from that of regression

- The quality assessor operates as if it is processing a single type of images

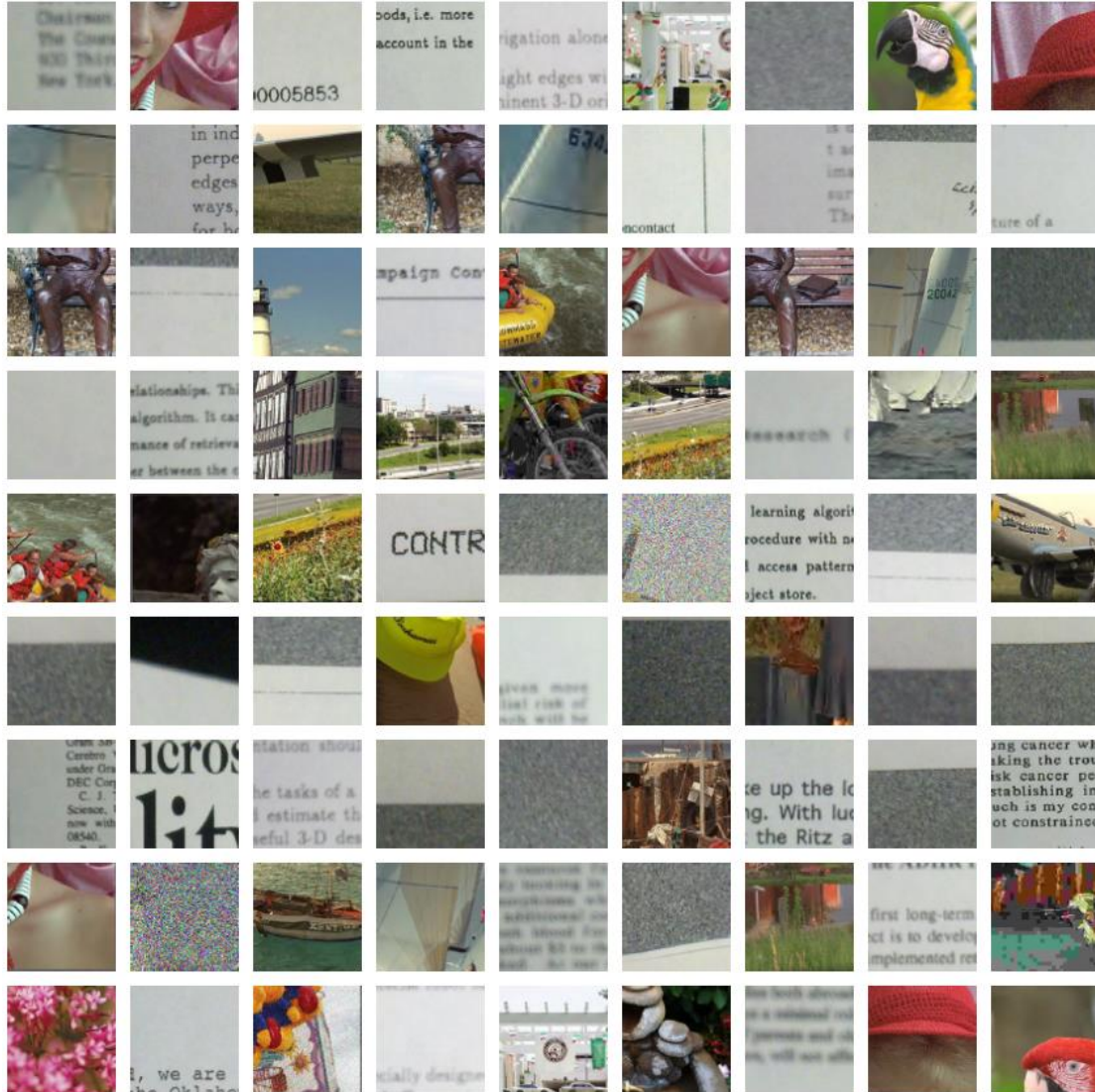# Unified Image Quality Assessment based on Contractive GAN



➤ **Main Objective:**

$$\min_{R,f} \max_{D} \left\{ \mathbb{E}_{x \sim p_B}\{\log[D(f(x))]\} + \mathbb{E}_{x \sim p_A}\{\log[1 - D(f(x))]\} \right.$$

$$\left. + \mathbb{E}_{x \sim p_A}\{|R(f(x)) - t_A|\} + \mathbb{E}_{x \sim p_B}\{|R(f(x)) - t_B|\} \right\}$$

➤ **Quality Discriminator:**

$$\max_{D_A}\left\{ \mathbb{E}_{x \sim p_A}\{\log[1 - |D_A(f(x)) - \sigma(x)|]\} \right\} +$$

$$\max_{D_B}\left\{ \mathbb{E}_{x \sim p_B}\{\log[1 - |D_B(f(x)) - \sigma(x)|]\} \right\}$$

where:

$$\sigma(x) = \begin{cases} 1, & \text{if } |R(x) - t| \leq \epsilon \\ 0, & \text{otherwise} \end{cases}$$

- Qualitative evaluation: visualization of the operation of the C-GAN model

# Unified Image Quality Assessment based on Transfer Learning

**LIVE + SOC**

**CSIQ + SOC**

# Unified Image Quality Assessment based on Contractive GAN

- Comparing to content-specific IQA and DIQA models

| IQA Models | CSIQ | | SOC | |
|---|---|---|---|---|
| | PLCC | SRCC | PLCC | SRCC |
| BLIINDS2 | - | 0.880 | N.A. | |
| DIQA | - | 0.870 | | |
| CORNIA | - | 0.854 | | |
| NRSL | - | **0.896** | | |
| CNN | N.A. | | 0.950 | 0.898 |
| CNN | | | 0.926 | 0.857 |
| RNN | | | **0.956** | 0.916 |
| LDA | | | - | 0.913 |
| Sparse Model | | | 0.935 | **0.928** |
| Proposed method | 0.92 | 0.89 | 0.932 | 0.916 |

- Cross-dataset evaluation of the proposed UIQA model on the natural scene images

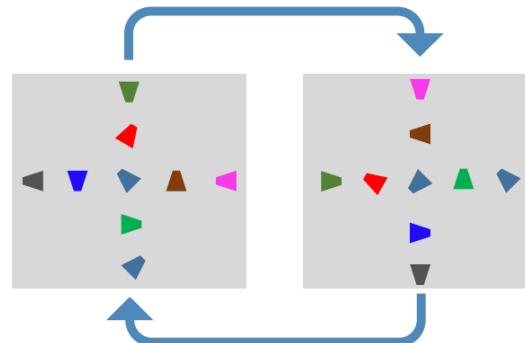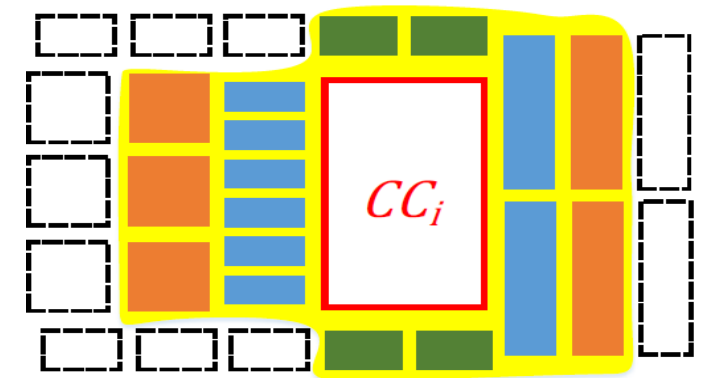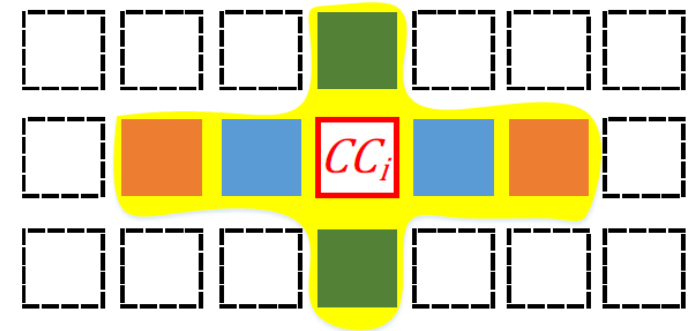| IQA Models | LIVE | |
|---|---|---|
| | PLCC | SRCC |
| BLIINDS2 | - | 0.915 |
| DIQA | - | **0.962** |
| CORNIA | - | 0.957 |
| NRSL | - | 0.808 |
| Proposed method | 0.91 | 0.952 |

➤ *Proximity*
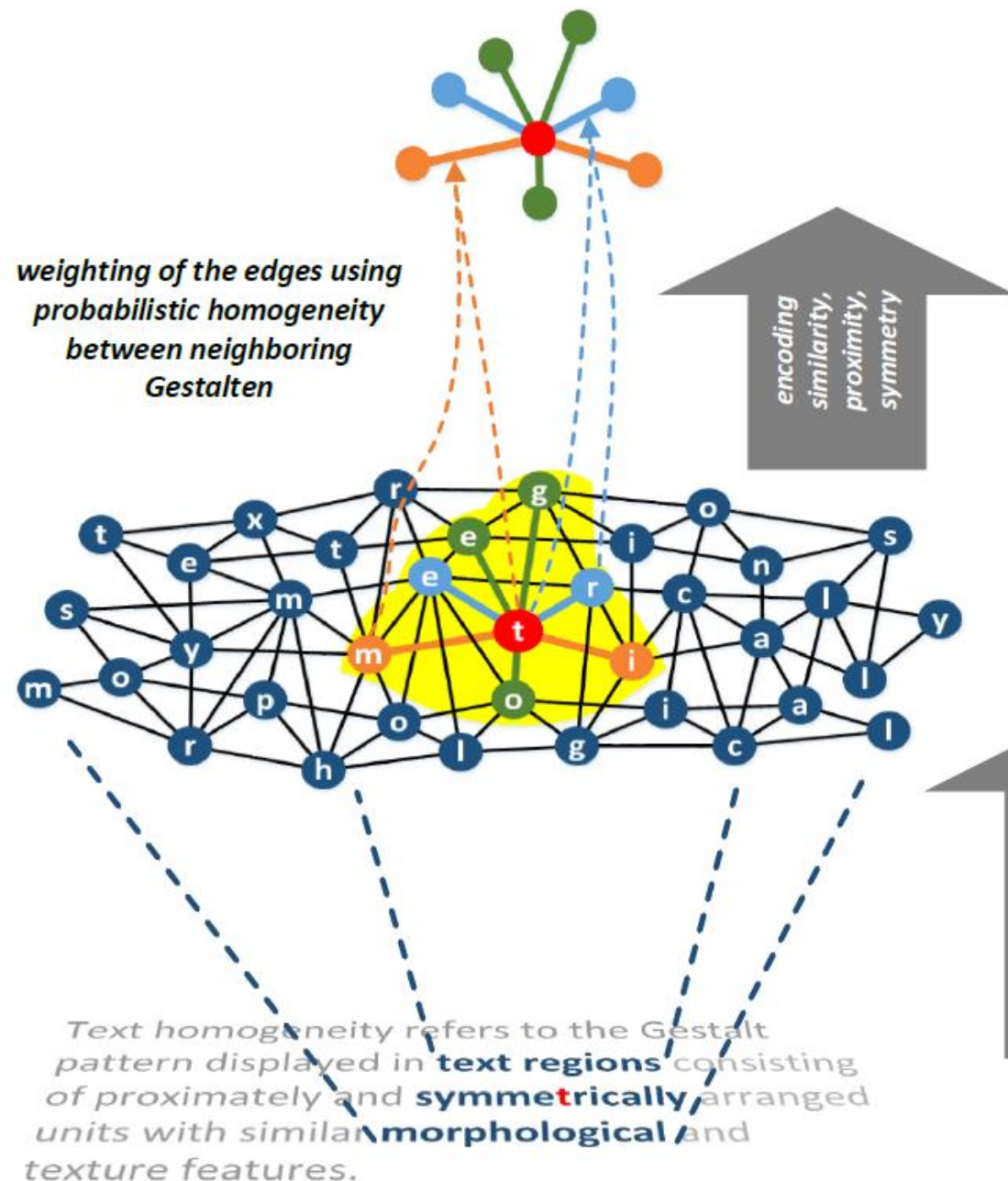
➤ *Similarity*

➤ *Symmetry*

➤ *Conceptualization*

**Text homogeneity** is the homogeneous pattern displayed in text regions, which consists of **proximately** and **symmetrically** arranged units with **similar** *morphological* and *texture* features.

$CC_i$

$CC_i$

weighting of the edges using
probabilistic homogeneity
between neighboring
Gestalten

encoding
similarity,
proximity,
symmetry

encoding geometric
configuration in the
neighborhoods

Text homogeneity refers to the Gestalt
pattern displayed in **text regions** consisting
of proximately and **symmetrically** arranged
units with similar **morphological** and
texture features.

➢ ***Description of local text homogeneity on G(V,E)***

If we take a one-step walk from a Gestalt $CC_i$ by following
an arbitrary (**symmetry**) direction, and arrives at another
Gestalt, say $CC_j$, the probability that $CC_j$ is located within
a short (**proximity**) distance and resembles (**similarity**) $CC_i$
is higher when $CC_i$ is a text component (e.g. a letter from
a paragraph).

• probabilistic weighting $w_{ij} = P(S_{ij} = s_{ij}^+)$

$$S_{ij} = \begin{cases} s_{ij}^+, & \text{if } CC_i \text{ and } CC_j \text{ are homogeneous,} \\ s_{ij}^-, & \text{if } CC_i \text{ and } CC_j \text{ are heterogeneous;} \end{cases}$$
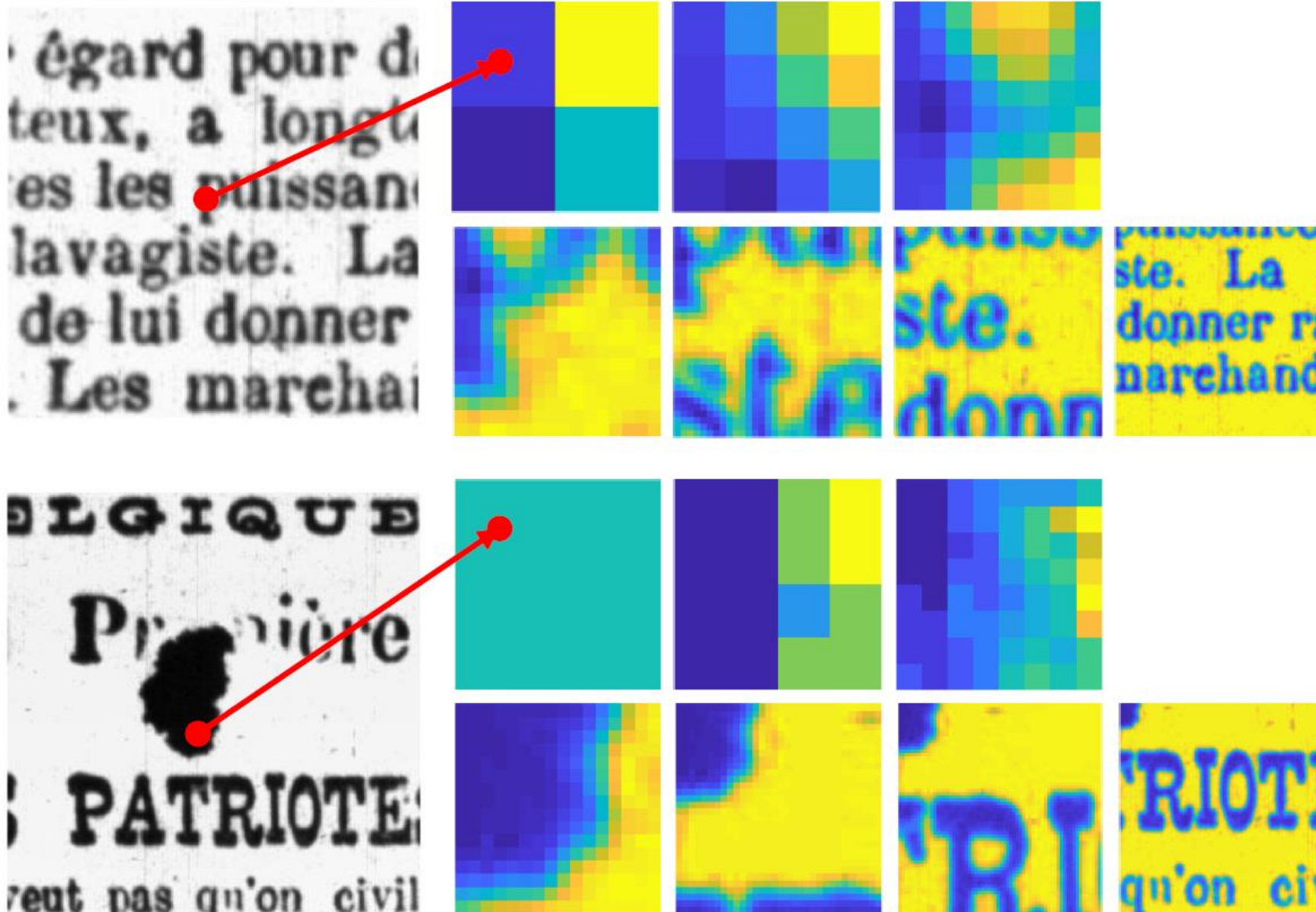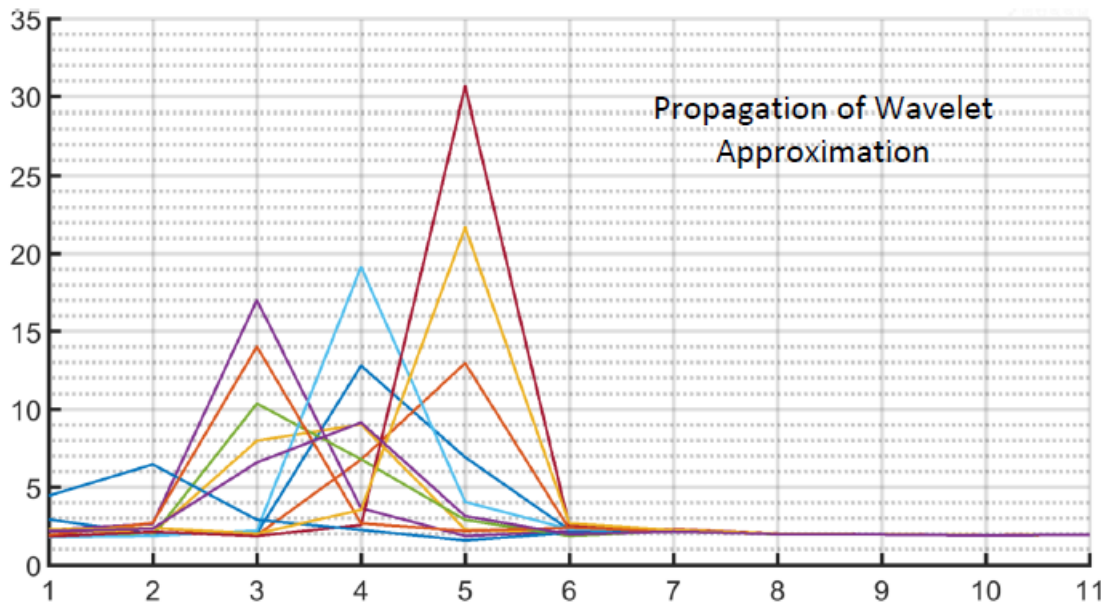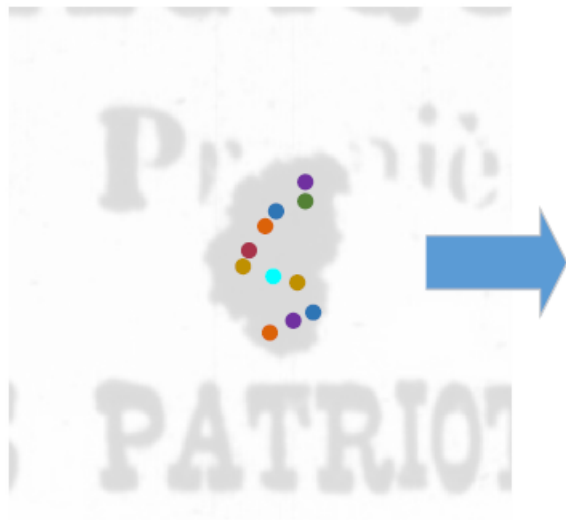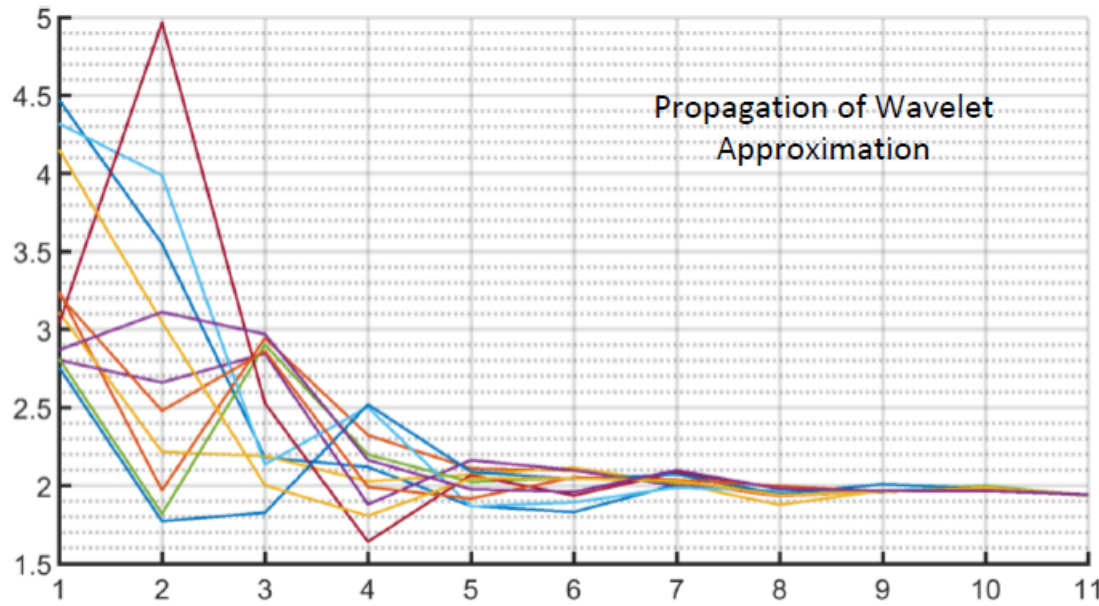
# Propagation of Wavelet Approximation



- **Text Homogeneity Revisit**
- Text Homogeneity pattern
- Neighborhood transition

# Propagation of Wavelet Approximation



> **Wavelet Propagation**

propagation of wavelet approximation (PWA) and propagation of cone-of-influence wavelet approximation (PCWA).
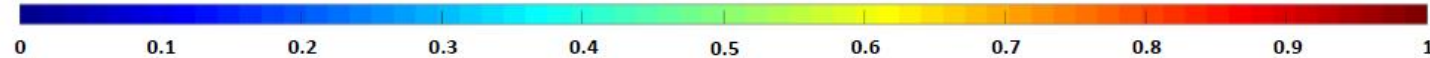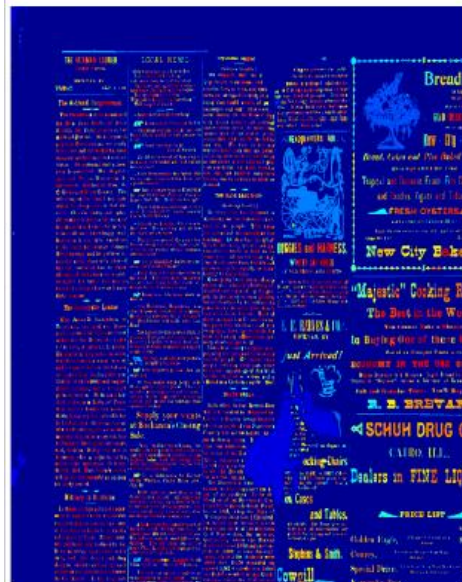
- PWA

$$\alpha_{n \to k, l} \triangleq \log_2 \left( \frac{1}{k-n} \sum_{j=n}^{k-1} \frac{|w_\phi^{j+1,l}|}{|w_\phi^{j,l}|} \right)$$
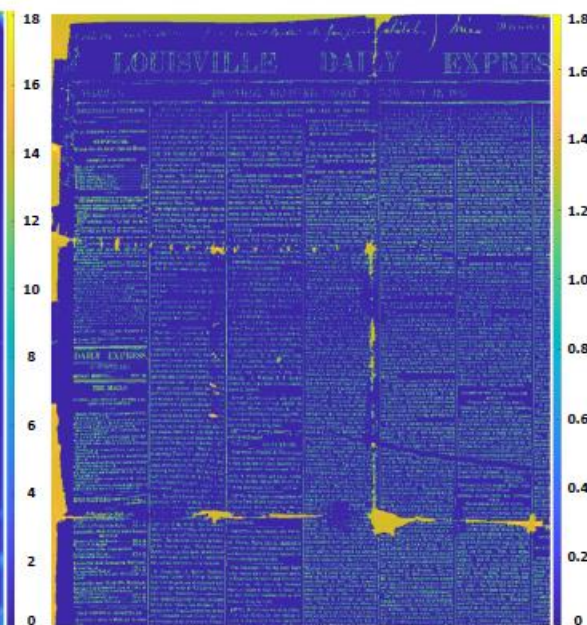
- PCWA
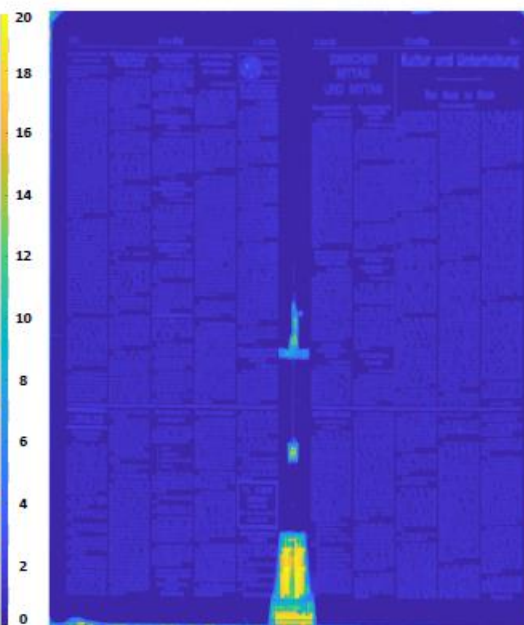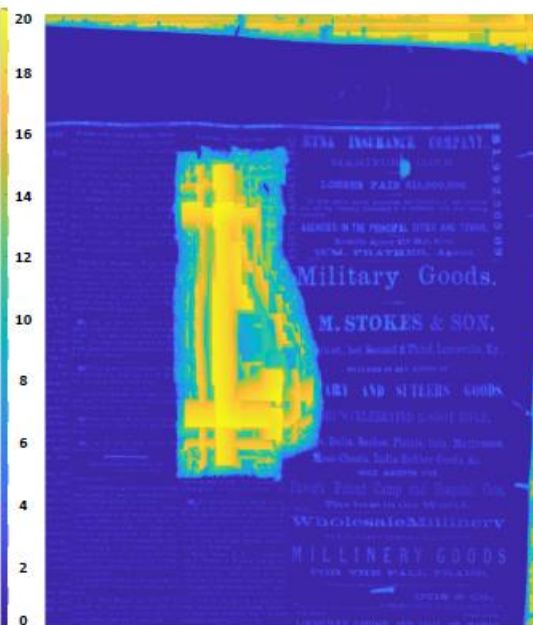
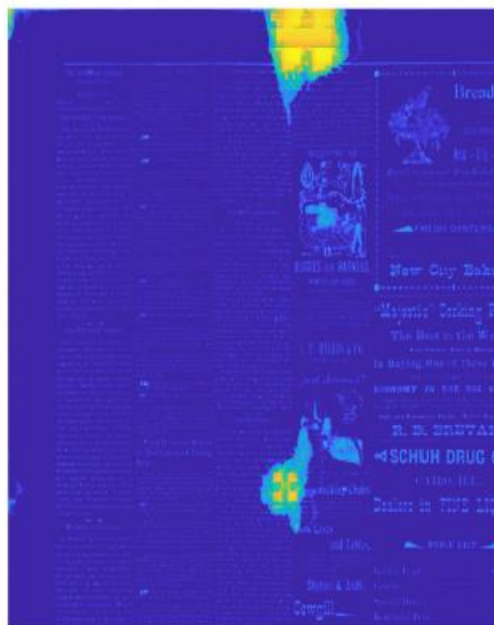$$\beta_{n \to k, l} \triangleq \log_2 \left( \frac{1}{k-n} \sum_{j=n}^{k-1} \frac{|I_{j+1,l}|}{|I_{j,l}|} \right),$$
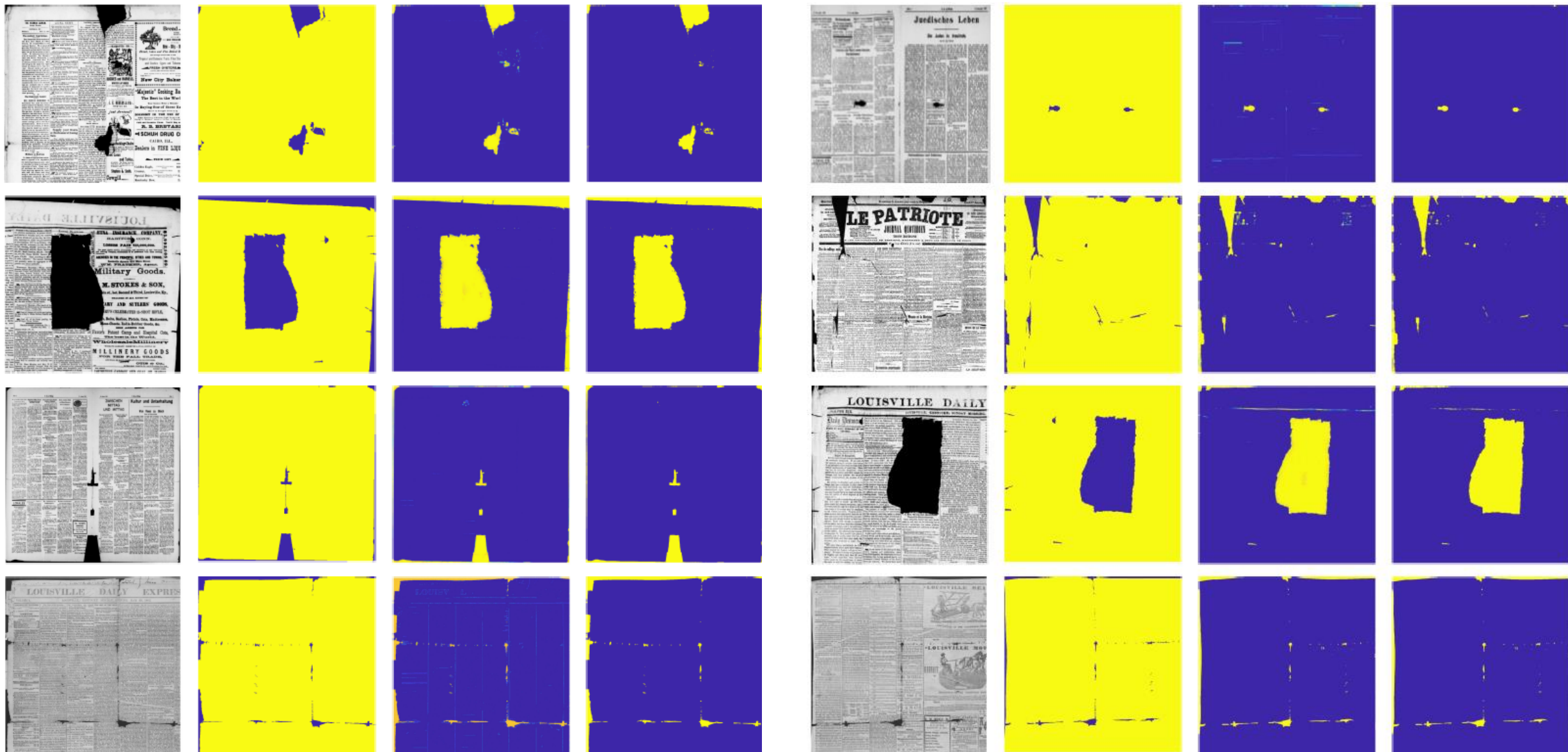
$$I_{j,l} \triangleq \sum_{m \in C(j,l)} |w_\phi^{j,m}|$$

# Bayesian Distortion Recognition

# Bayesian Distortion Recognition

# Bayesian Distortion Recognition

Thank you for your attention!

Q & A