# Deterministic Routing between Layout Abstractions for Multi-Scale Classification of Visually Rich Documents

Ritesh Sarkhel*, Arnab Nandi
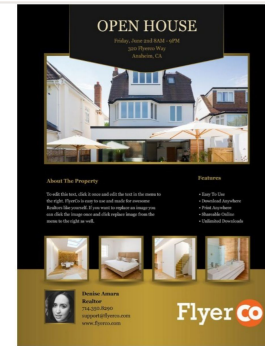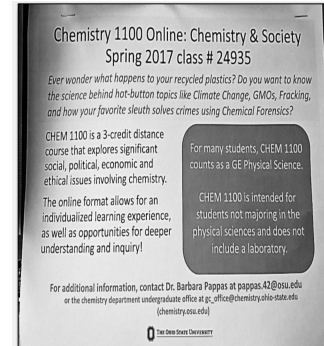Department of Computer Science and Engineering

THE OHIO STATE UNIVERSITY

https://interact.osu.edu/

# Outline

❖ **Visually Rich Documents**

❖ **Challenges**

❖ **Our contributions**

  ➢ **Spatial Pyramid Model**

  ➢ **Deterministic Routing Scheme**

❖ **End-to-end architecture for classification**

❖ **Experimental results**

  ➢ **Tobacco-Litigation dataset**

  ➢ **Outperforms previous state-of-the-art by 4.73%**

THE OHIO STATE UNIVERSITY

https://interact.osu.edu/

# Visually Rich Document

- Visual cues to highlight distinct semantic entities
  - Color
  - Font size
  - Text orientation and positioning
  - Symmetry etc.

- Classification is a precursor of many document understanding tasks
  - Indexing
  - Information extraction

https://interact.osu.edu/

# Our objective

1. **Maximize classification accuracy using visual (dis)similarity in layouts**
   - Multi-scale features
     - Encode local invariant patterns from multiple resolutions of the document
     - More robust than single-scale features: Leads to better classification accuracy

2. **Minimize end-to-end inference turnaround**
   - Suitable for modern interactive workflows

Assumption: A soft constraint on the relative positioning of components in documents belonging to the same class
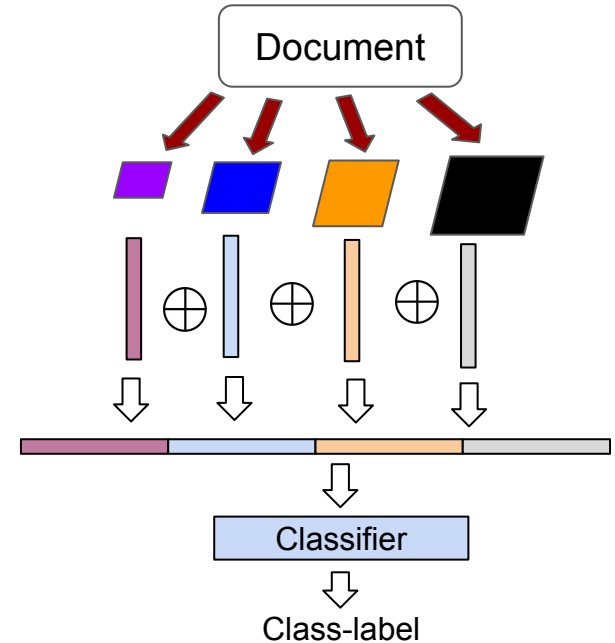
https://interact.osu.edu/

# Outline

❖ Visually Rich Documents

❖ **Challenges**

❖ Our contributions

➢ **Spatial Pyramid Model**

➢ **Deterministic Routing Scheme**

❖ End-to-end architecture for classification

❖ Experimental results

➢ **Tobacco-Litigation dataset**

➢ **Outperforms previous state-of-the-art by 4.73%**

THE OHIO STATE UNIVERSITY

https://interact.osu.edu/

# Challenges

- ## Multi-scale: Aggregation based methods
    - Increased overhead in inference turnaround
    - Marginal gain in accuracy from single-scale counterparts

- ## Infeasible when bound by near real-time latency
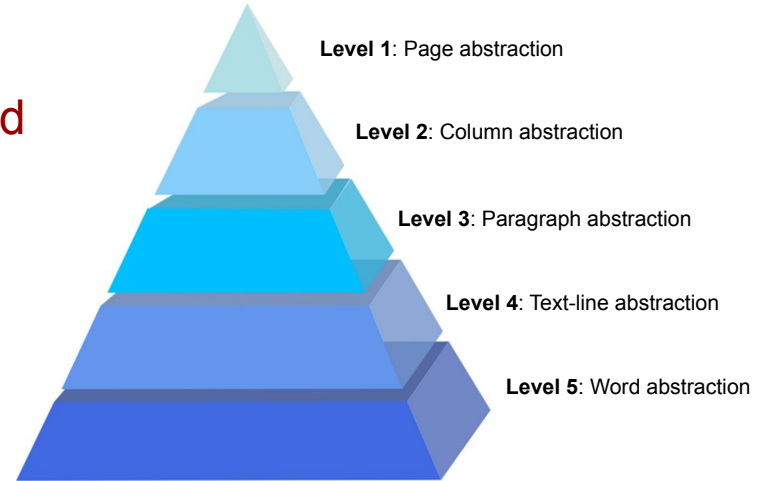    - e.g. Interactive workflow on edge-devices

https://interact.osu.edu/

# Outline

❖ Visually Rich Documents

❖ Challenges and observations

❖ **Our contributions**

➢ **Spatial Pyramid Model**

➢ **Deterministic Routing Scheme**

❖ Neural architecture for classification

❖ Experimental results

➢ **Tobacco-Litigation dataset**

➢ **Outperforms previous state-of-the-art by 4.73%**

THE OHIO STATE UNIVERSITY

https://interact.osu.edu/

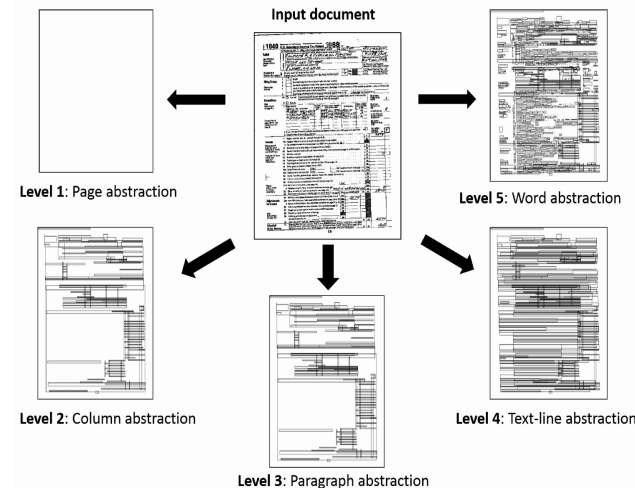# A Spatial Pyramid Model for document representation

- Objective: Increase accuracy by extracting highly discriminative features

- Represent the document as a spatial pyramid
  - Each level corresponds to an abstraction in layout hierarchy:
    - Words
    - Text-line
    - Paragraph
    - Column
    - Page



**Level 1**: Page abstraction

**Level 2**: Column abstraction

**Level 3**: Paragraph abstraction

**Level 4**: Text-line abstraction

**Level 5**: Word abstraction

https://interact.osu.edu/

# Layout abstraction at each level of the Spatial Pyramid Model

- **Each pyramid level corresponds to an abstraction at that level of the layout hierarchy**
  - Binary Image with same dimensions as the input document
  - e.g. Word abstraction
    - Bounding box corresponding to each word of the document

- **Layout hierarchy derived using an open-source page segmentation algorithm [1,2]**



Input document

Level 1: Page abstraction

Level 5: Word abstraction

Level 2: Column abstraction

Level 4: Text-line abstraction

Level 3: Paragraph abstraction

[1] https://github.com/tesseract-ocr/tesseract
[2] Smith, Ray. "An overview of the Tesseract OCR engine." In *Ninth International Conference on Document Analysis and Recognition (ICDAR 2007)*, vol. 2, pp. 629-633. IEEE, 2007.
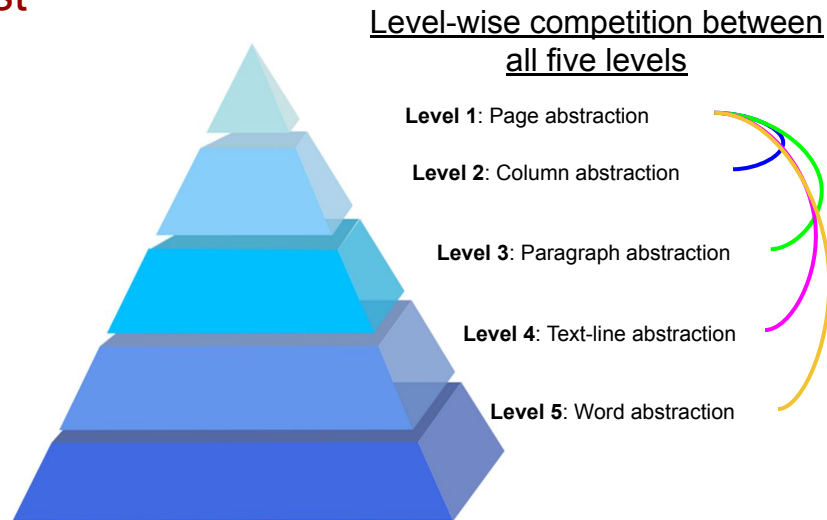
THE OHIO STATE UNIVERSITY

https://interact.osu.edu/

# Outline

❖ Visually Rich Documents

❖ Challenges and observations

❖ **Our contributions**

    ➢ **Spatial Pyramid Model**

    ➢ **Deterministic Routing Scheme**

❖ Neural architecture for classification

❖ Experimental results

    ➢ **Tobacco-Litigation dataset**

    ➢ **Outperforms previous state-of-the-art by 4.73%**

THE OHIO STATE UNIVERSITY

https://interact.osu.edu/

# A Deterministic Routing Scheme for feature selection

- Attention-like operation to select the most discriminative pyramid level

- Introduce level-wise competition
  - Select the most discriminating pyramid level

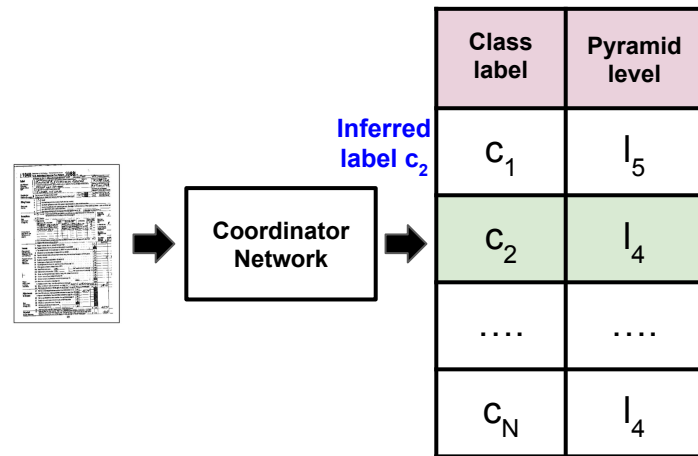- The winner represents the document for final prediction

Level-wise competition between all five levels

**Level 1**: Page abstraction

**Level 2**: Column abstraction

**Level 3**: Paragraph abstraction

**Level 4**: Text-line abstraction

**Level 5**: Word abstraction

https://interact.osu.edu/

# Supervised coordinator network for Deterministic Routing

- ## LadderNet: Softmax classifier network
  - Extended the MobileNetV2 architecture [1,2]
  - Fast: Factorized convolution
  - Pretrained on ImageNet, fine-tuned on training corpus

- ## Softmax label to query the multiplex table
  - Single entry for each class-label
  - Route to the corresponding pyramid level

| Class label | Pyramid level |
|:---:|:---:|
| $c_1$ | $l_5$ |
| $c_2$ | $l_4$ |
| …. | …. |
| $c_N$ | $l_4$ |

**Inferred label $c_2$**

**Coordinator Network**

[1] Howard, Andrew G., Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. "Mobilenets: Efficient convolutional neural networks for mobile vision applications." *arXiv preprint arXiv:1704.04861* (2017).

[2] Sandler, Mark, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. "Mobilenetv2: Inverted residuals and linear bottlenecks." In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4510-4520. 2018.
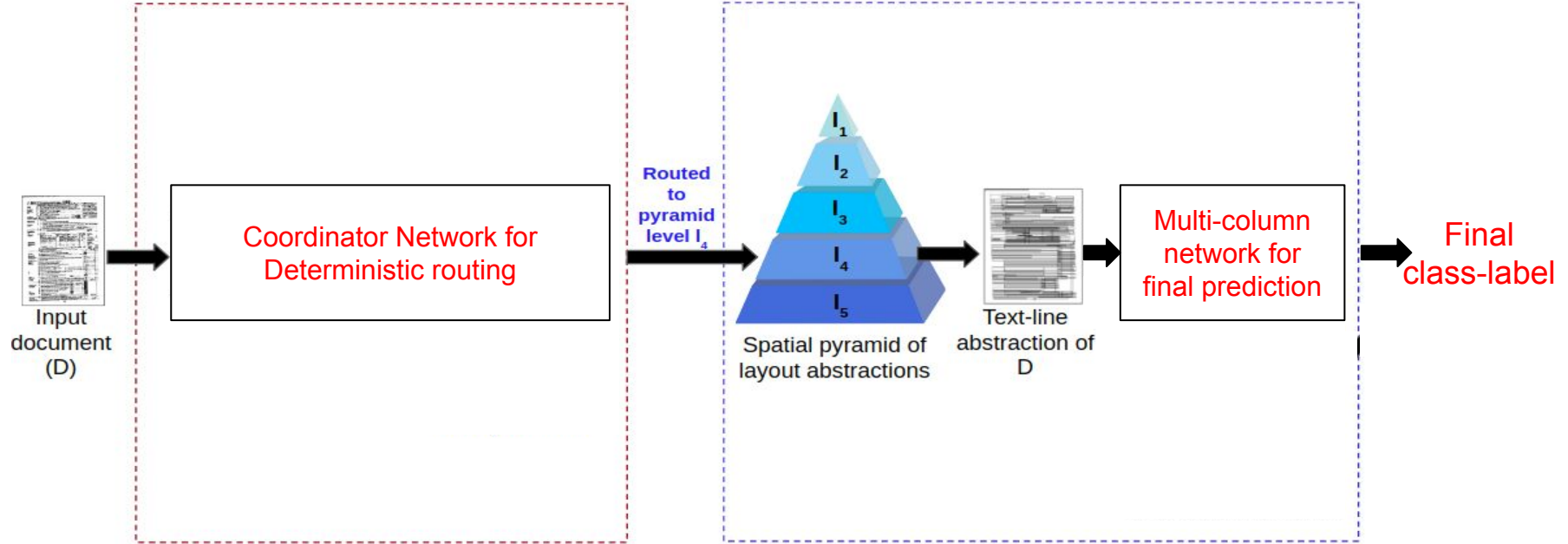
THE OHIO STATE UNIVERSITY

https://interact.osu.edu/

# Outline

❖ Visually Rich Documents

❖ Challenges and observations

❖ Our contributions

   ➢ **Spatial Pyramid Model**

   ➢ **Deterministic Routing Scheme**

❖ **End-to-end architecture for classification**

❖ Experimental results

   ➢ **Tobacco-Litigation dataset**

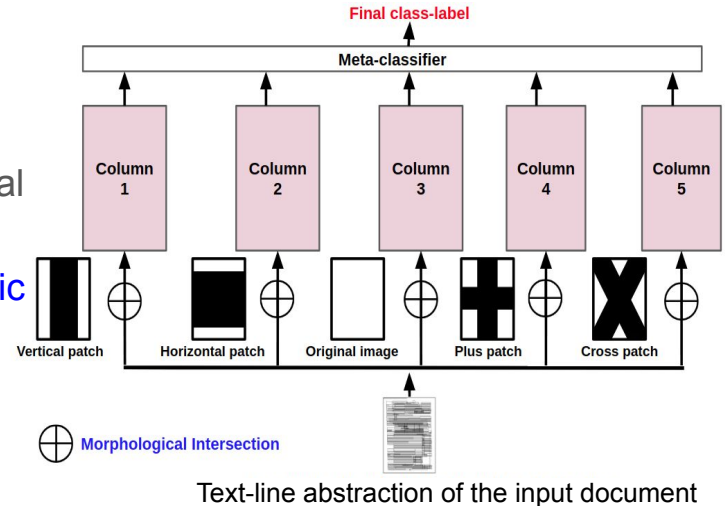   ➢ **Outperforms previous state-of-the-art by 4.73%**

THE OHIO STATE UNIVERSITY

https://interact.osu.edu/

# End-to-end classification workflow



Input document (D) → Coordinator Network for Deterministic routing → Routed to pyramid level $I_4$ → Spatial pyramid of layout abstractions ($I_1$, $I_2$, $I_3$, $I_4$, $I_5$) → Text-line abstraction of D → Multi-column network for final prediction → Final class-label

https://interact.osu.edu/

# A Multi-column Network for Final prediction

- ## Multi-column neural architecture
  - Depth-separable convolutional columns
  - <u>Input</u>: Morphological transformation of the document
  - Combine predictions from each column to compute final class label
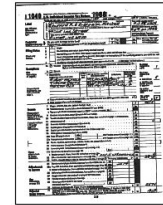    - Meta-classifier: Simple-avg, weighted-avg, logistic regression and MLP



Text-line abstraction of the input document

# Outline

❖ Visually Rich Documents

❖ Challenges and observations

❖ Our contributions

    ➢ **Spatial Pyramid Model**

    ➢ **Deterministic Routing Scheme**

❖ End-to-end architecture for classification

❖ Experimental results

    ➢ **Tobacco-Litigation dataset**

    ➢ **Outperforms previous state-of-the-art by 4.73%**

THE OHIO STATE UNIVERSITY

https://interact.osu.edu/

# Experiments

- **Four publicly available datasets of single-page documents**
  - NIST Special dataset-6 (easy)
  - Medical Article Records Groundtruth or MARG dataset (medium)
  - Tobacco Litigation dataset (hard)
  - IIT-CDIP dataset (hard)

- **Evaluation metrics**
  - Classification accuracy
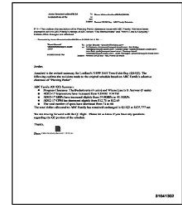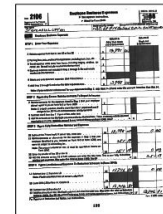  - Inference turnaround time
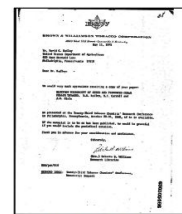


(A1) NIST    (B1) MARG    (C1) Tobacco    (D1) CDIP

(A2) NIST    (B2) MARG    (C2) Tobacco    (D2) CDIP

https://interact.osu.edu/

# Result highlights

- **Tobacco Litigation dataset**
  - High inter-class and low intra-class layout similarity
  - 3482 documents, 10 document classes (e.g. 'letter', 'memo', 'email' etc.)

- **Median classification accuracy over 25 trials := 82.78%**
  - Absolute improvement over previous state-of-the-art := 4.73%

- **Average inference turnaround time := 362 (±10.27) ms**
  - Average speed-up factor > 6

**More results and ablation study in paper !!**

THE OHIO STATE UNIVERSITY

https://interact.osu.edu/

# Takeaways

- **Multi-scale classifier for visually rich documents**
  - Spatial pyramid model
    - Highly discriminative features → Increase in accuracy
  - Level-wise competition for optimal document representation
    - Attending to the most discriminative layout abstraction improves accuracy

- **Inference turnaround time < interactive latency (≅ 500 ms)**

- **Robust performance on four publicly available datasets**

THE OHIO STATE UNIVERSITY

https://interact.osu.edu/

# Thank you!