

第四届全国文字与计算学术研讨会，2014.10.25-26，北京

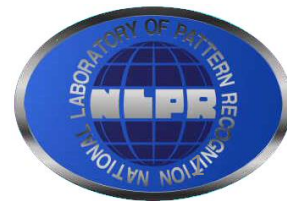
文档分析技术研究现状与趋势

刘成林

中国科学院自动化研究所
模式识别国家重点实验室

liucl@nlpr.ia.ac.cn

<http://www.nlpr.ia.ac.cn/liucl>



Outline

- 文档与文档分析
- 文档分析研究问题
- 领域发展简史
- 研究现状
 - 主要方法
 - 性能状况
- 国内现状
- 趋势与展望

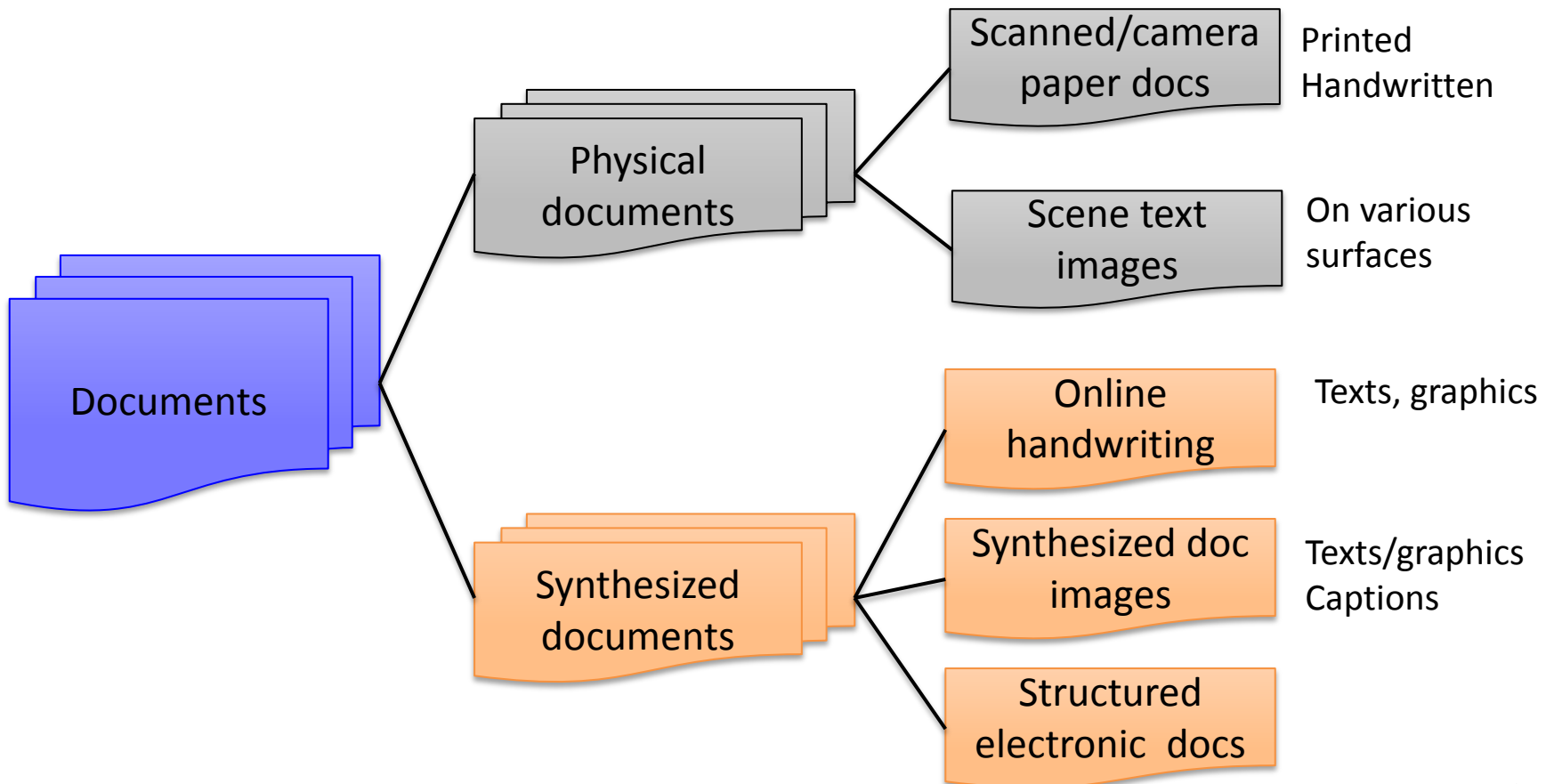
文字识别与文档分析

- 文字识别(Character Recognition)
 - 字符图像转换为符号代码
- 文档分析(Document Analysis)
 - 从文档图像提取文本信息
 - 包括文本分割、识别、上下文处理、语义信息提取等
- 文档分析的意义
 - 数据压缩
 - 内容理解/语义提取
- 与文本分析(Text Analysis, NLP)的关系
 - Text Analysis: 从电子文本开始
 - DA: 从图像开始
 - 电子文档(如PDF): boundary, DA利用其结构信息



文档的种类

- 什么是文档
 - 载有文字符号的纸张、图像或电子文件

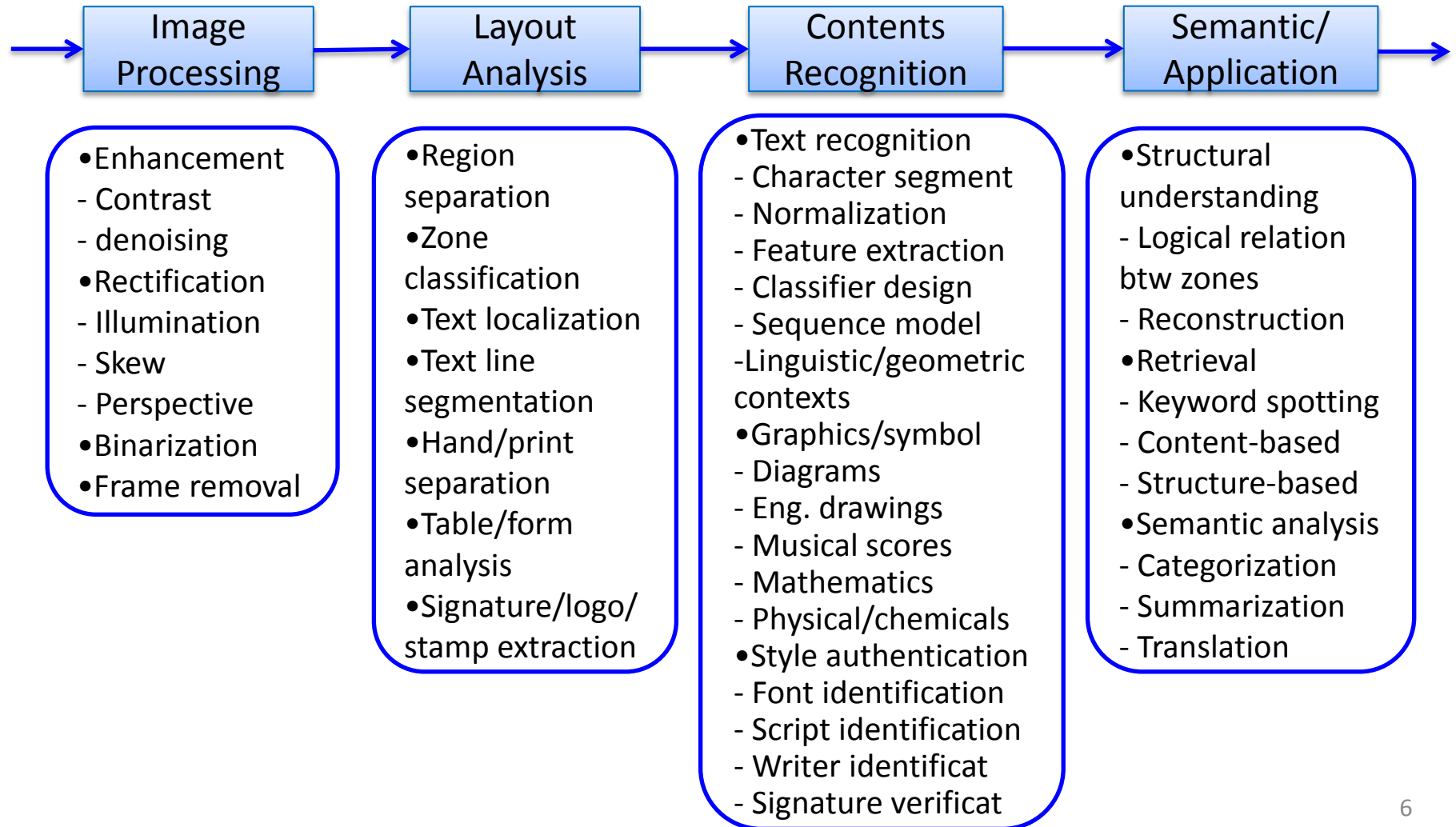


具体种类与应用

- Paper Documents
 - Books, journals, newspapers, letter/parcel envelopes, certificates, notes, forms, business cards, engineering drawings, musical scores
- Scene Texts
 - Sign boards, license plate, street numbers
 - Texts on wood, metal, cloth, oracle bones, etc.
- Online Handwriting
 - Texts, graphics, signature, mathematics, sketch, gesture
- Synthesized Document Images
 - Web doc images
 - Captions on image/video
- Structured Documents
 - Web pages
 - Word doc, PDF, RTF, etc.

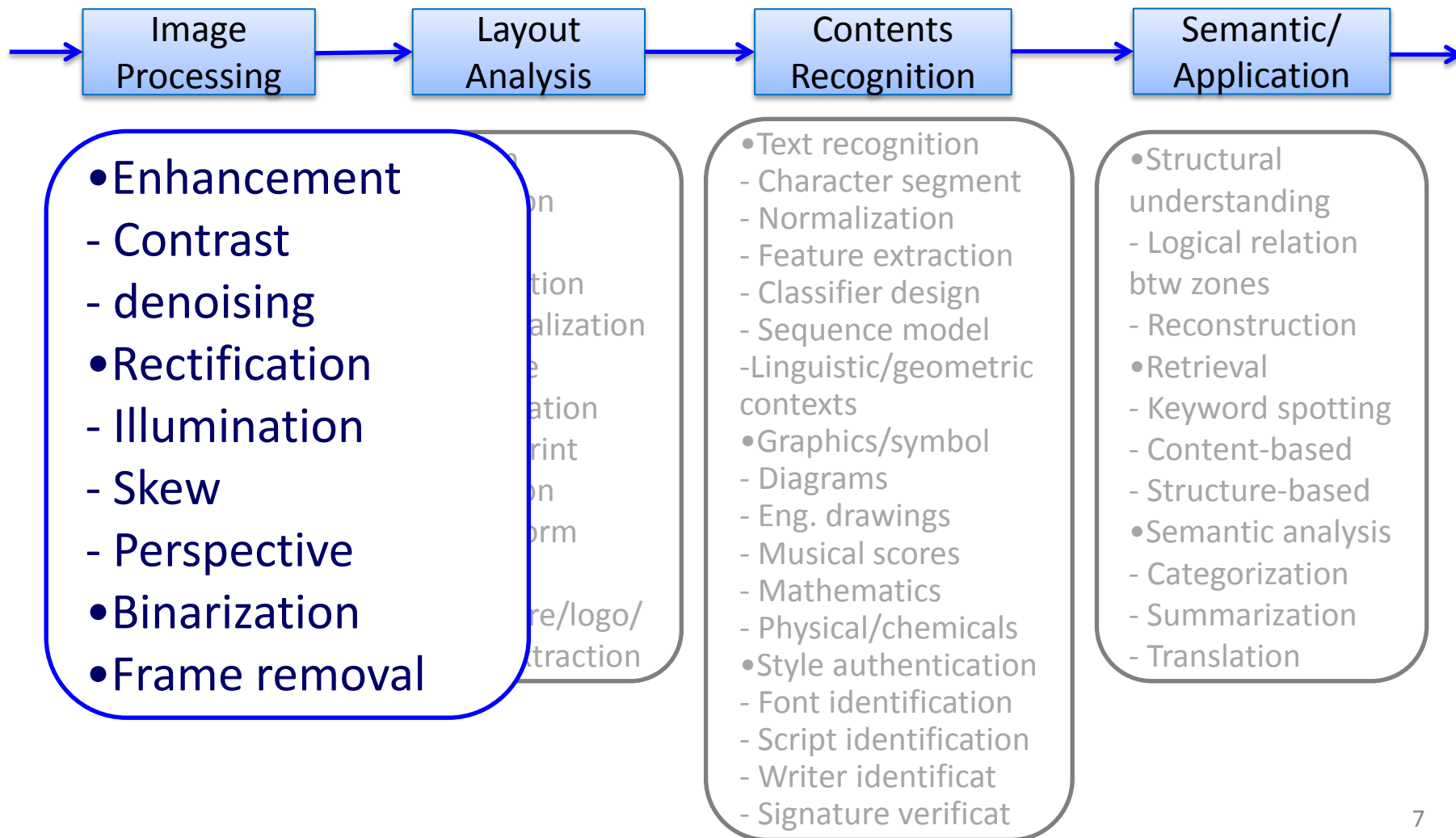
文档分析研究问题

- From Image to Semantics



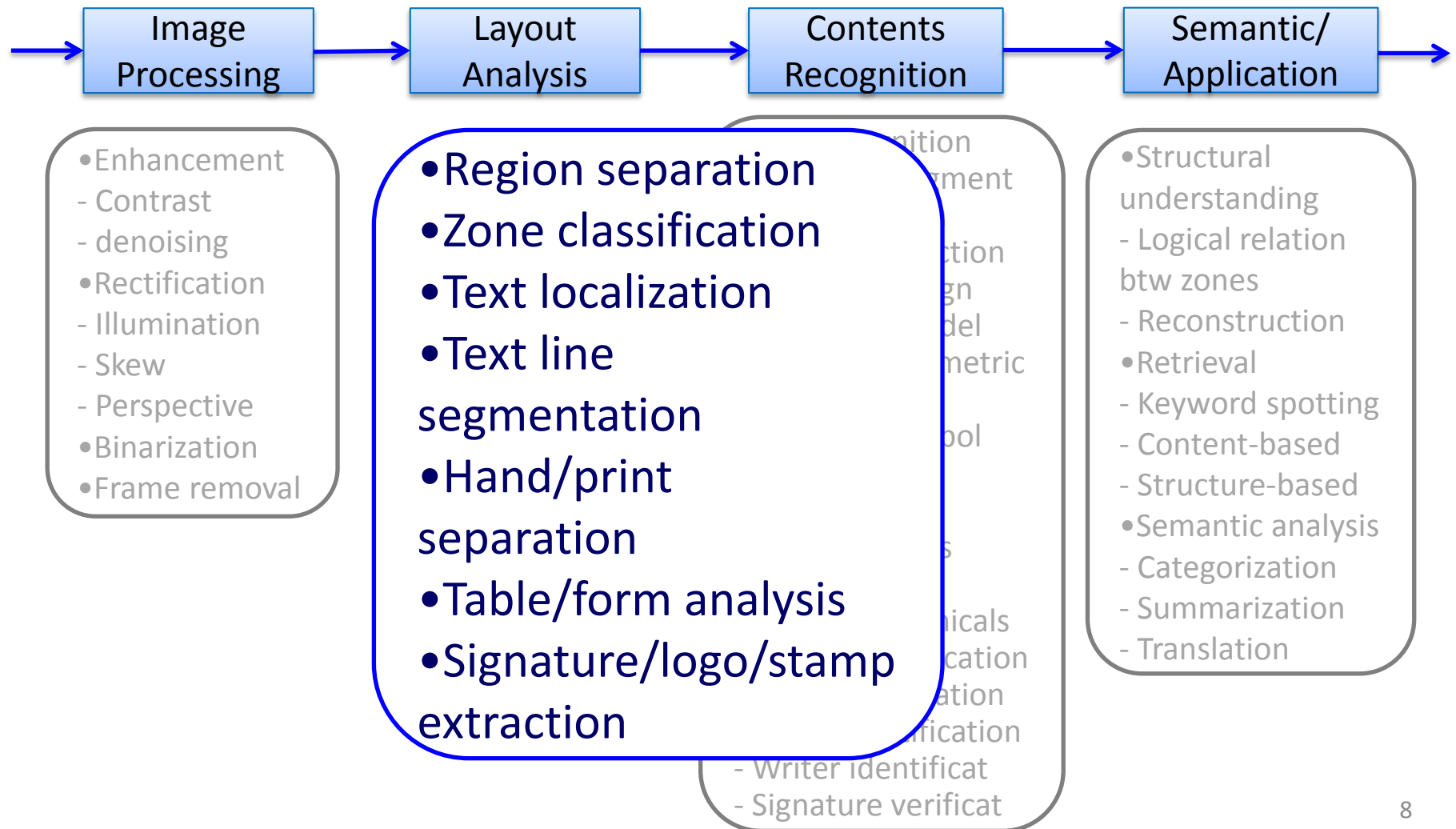
文档分析研究问题

- From Image to Semantics



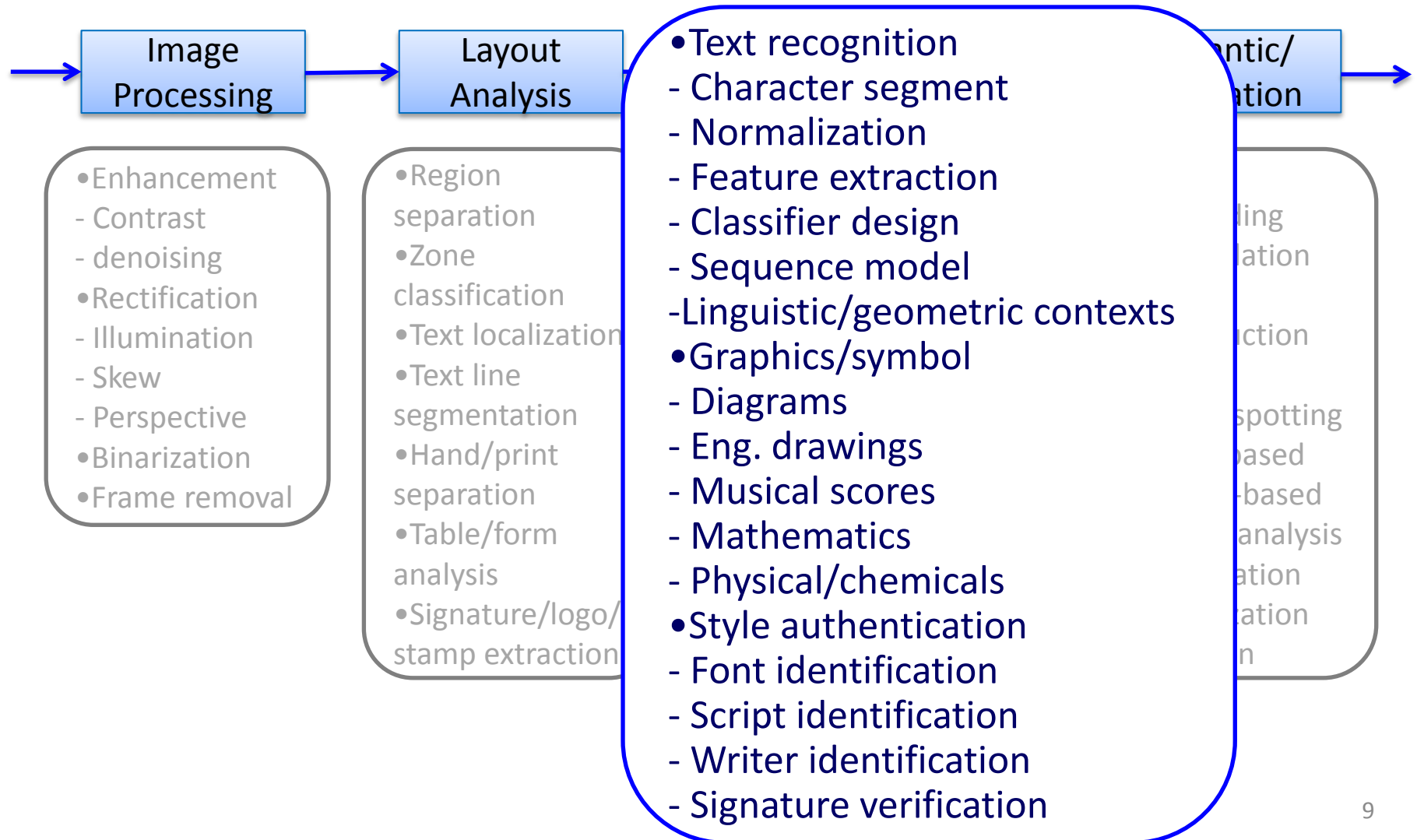
文档分析研究问题

- From Image to Semantics



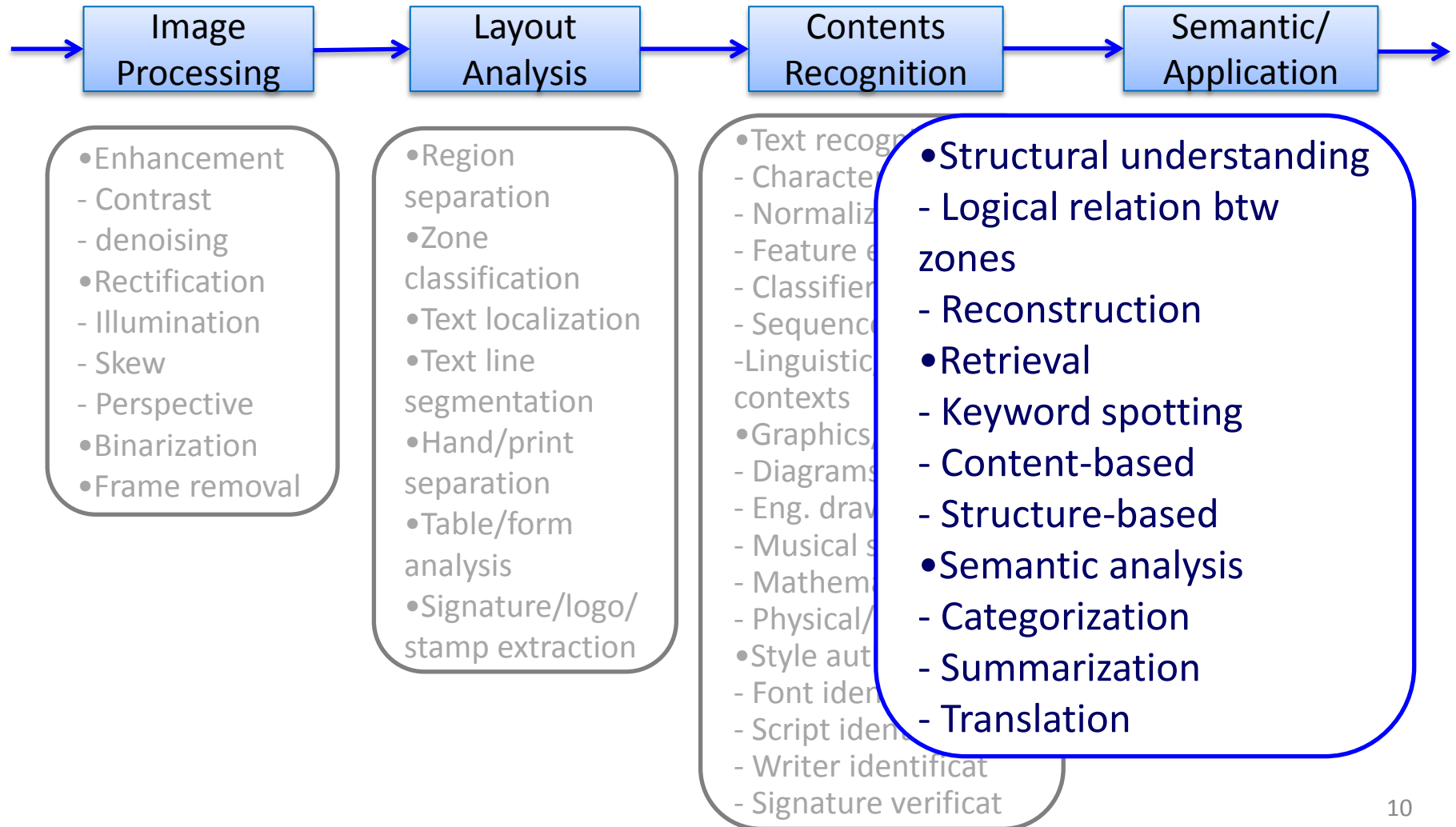
文档分析研究问题

- From Image to Semantics



文档分析研究问题

- From Image to Semantics

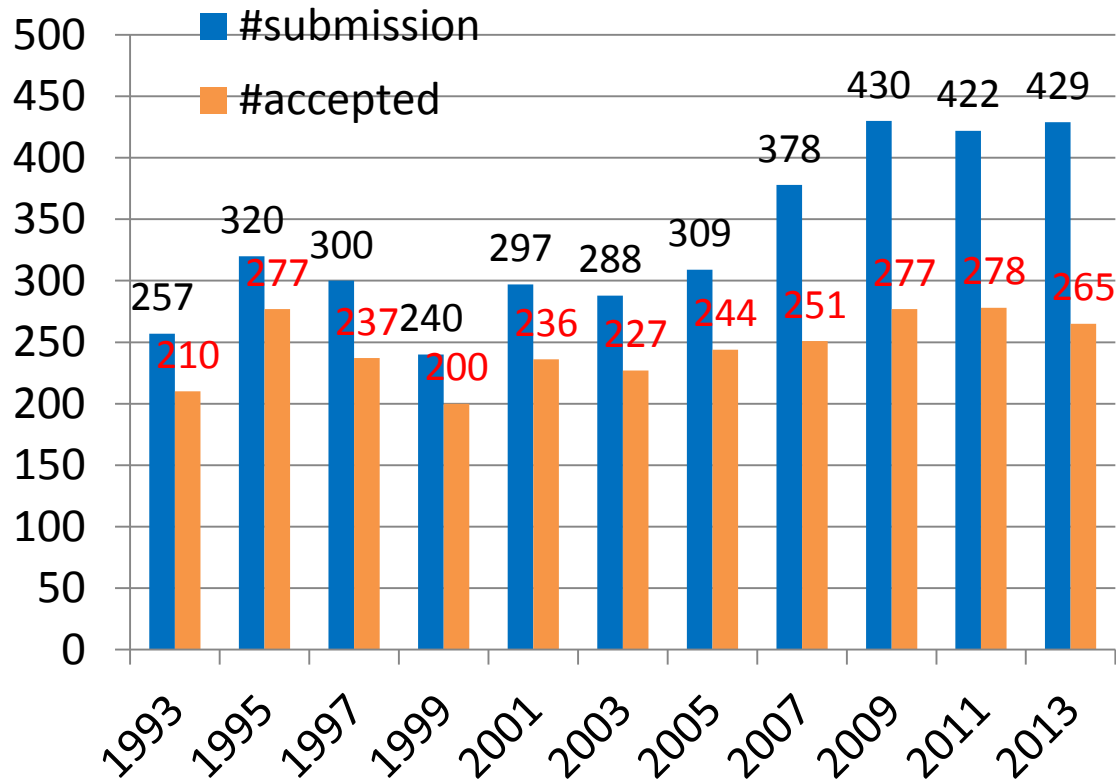


领域发展简史

Time	Methods	Target/Application	Events
1920s	Optical template matching	Printed digits/letters	1 st patent on OCR
1950s-1960s	Correlated matching, simple structural analysis	Printed digits/letters Printed Chinese (1966)	1 st PR Workshop in 1966
1970s-1980s	Feature matching (normalization, feature extraction), Structural matching, Statistical PR	Handprinted digits/letters Printed/handprinted words Handprinted Japanese/Chinese	1 st ICPR in 1972 IAPR founded in 1978
1990s	Research of various issues, including layout analysis and segmentation	Practical applications in various areas (document entry, mail sorting, forms, business cards, text input)	PC got popular Internet 1 st IWFHR/ICDAR/DAS in 1990/91/94
2000s	Re-inventing existing methods (e.g., HMM, NN) Borrowing from ML and CV (e.g., SSL, BoW)	Remaining hard problem Improve existing apps Explore new apps (e.g., camera-based, ink documents)	Google, Baidu Facebook, twitter Weibo, smart phone Mobile Internet

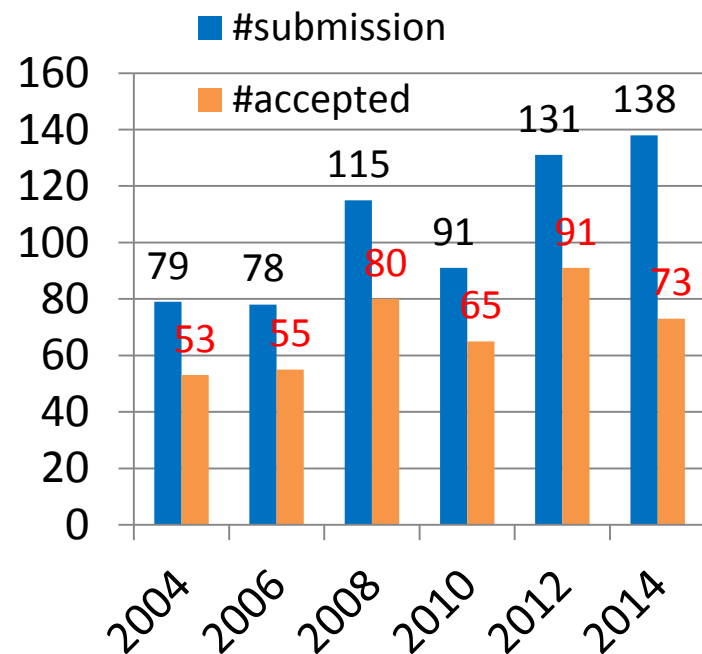
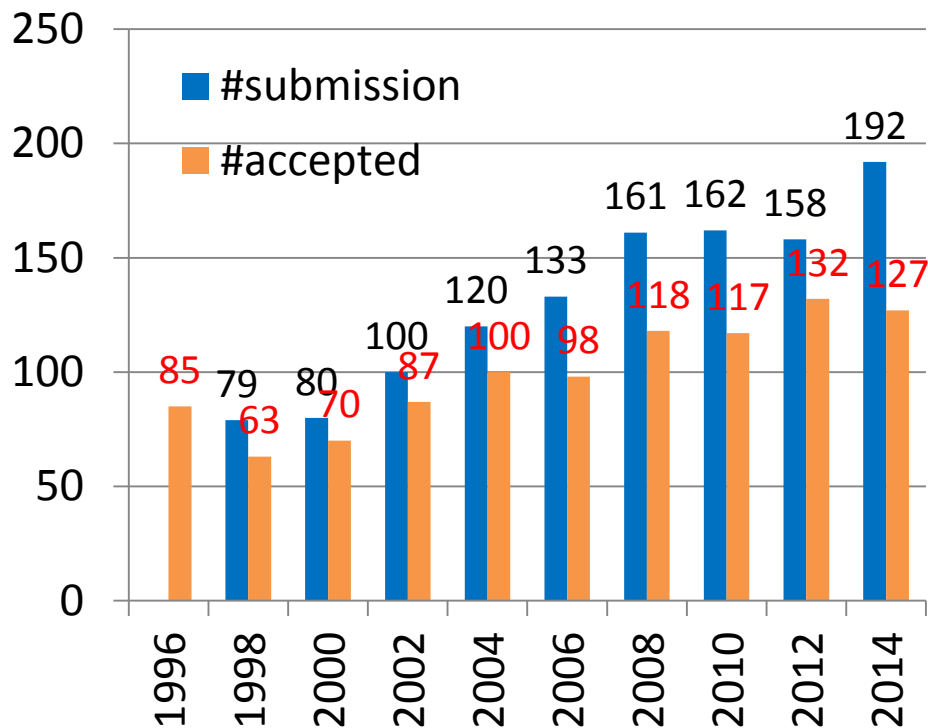
会议规模

- ICDAR: Int'l Conf. on Document Analysis and Recognition (bi-ennial from 1991)



会议规模和主要话题

- ICFHR: Int'l Conf. on Frontiers of Handwriting Recognition (from 1990, formerly IWFHR)
- DAS: IAPR Int'l Workshop on Document Analysis Systems (bi-ennial from 1994)



主要话题(Oral Sessions)

- ICDAR2013

- Character recognition 2
- Applications 2
- Document image processing 2
- Online handwriting recognition
- Binarization
- Handwriting recognition
- Scene text segmentation
- Keyword spotting
- Video and camera-based OCR
- Document analysis and classification
- Handwritten text recognition 2
- Layout analysis
- Signature verification
- Graphics recognition 2
- Historical documents

- 2011

- Datasets and performance evaluation
- Mathematics recognition
- Document segmentation
- Forensic document analysis
- Document retrieval

- ICDAR1995

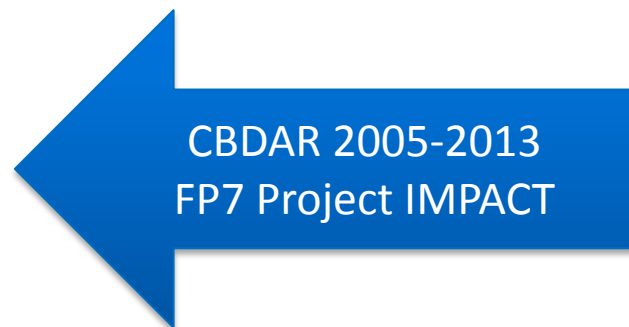
- Word recognition 2
- Neural networks 2
- Online recognition 2
- Application systems
- Handwritten character recognition
- Segmentation 2
- Image processing 2
- Database and document retrieval
- Signature verification
- Model based document analysis
- Map interpretation
- Classification methods
- Document understanding
- Feature extraction
- Reconstruction and interpretation
- Document analysis
- Symbol recognition
- Postprocessing
- Document structure and analysis
- Modeling methods
- Drawing and map reconstruction
- Shape models
- Item extraction
- Document segmentation
- Diagram/drawing understanding
- Word and character recognition
- Theoretical approach and music recognition
- Learning
- Layout analysis
- Skew detection and processing
- Storage and retrieval

主要话题(Oral Sessions)

- ICFHR2014
 - Word spotting
 - Document image pre-processing
 - Signature verification
 - Applications 2
 - Neural Networks for Handwriting Recognition
 - Online handwriting recognition
 - Language models for handwriting recognition
 - HMMs for handwriting recognition
 - Writer identification
- 2012
 - Character recognition
 - Flowchart recognition
 - Word segmentation
 - Multilingual recognition
 - Bank check processing and postal automation
 - Mathematical expression
 - Forensic applications
 - Historical documents
- IWFHR1998
 - Online handwriting recognition 2
 - Handwritten form processing
 - Handwritten word recognition
 - Segmentation
 - Oriental script processing
 - Numeral recognition
 - Emerging techniques
- 2002
 - Classifier design
 - Online recognition
 - Word recognition
 - Learning methods
 - Multiple classifiers
 - Pen computing and document applications
 - Signature verification and writer identification
 - Offline recognition

话题主要变化

- 减少的话题
 - Map and engineering drawings
 - Character recognition: no longer categorized into digit, oriental, feature, classifier and so on
 - Form and postal, although still challenging
- 增加的话题
 - Scene text
 - Camera and video-based
 - Historical documents
- 重点发生变化的话题
 - Layout analysis: from scanned printed documents to handwritten documents and camera-based documents
 - Retrieval: from printed to handwritten/historical, keyword spotting getting attention
 - Handwriting: from word to text (including language model)
 - Neural networks: from shallow to deep



研究现状：主要方法

- 图像预处理
 - Enhancement/denoising
 - Morphology
 - MRF
 - Deblurring
 - Binarization
 - Local/adaptive binarization
 - Stroke edges
 - Classification-based, MRF, CRF
 - Rectification
 - 3D shape modeling
 - Cylindrical surface reconstruction
 - Polynomial curve fitting with text lines
 - Frame line removal
 - Graph-based, heuristic rules

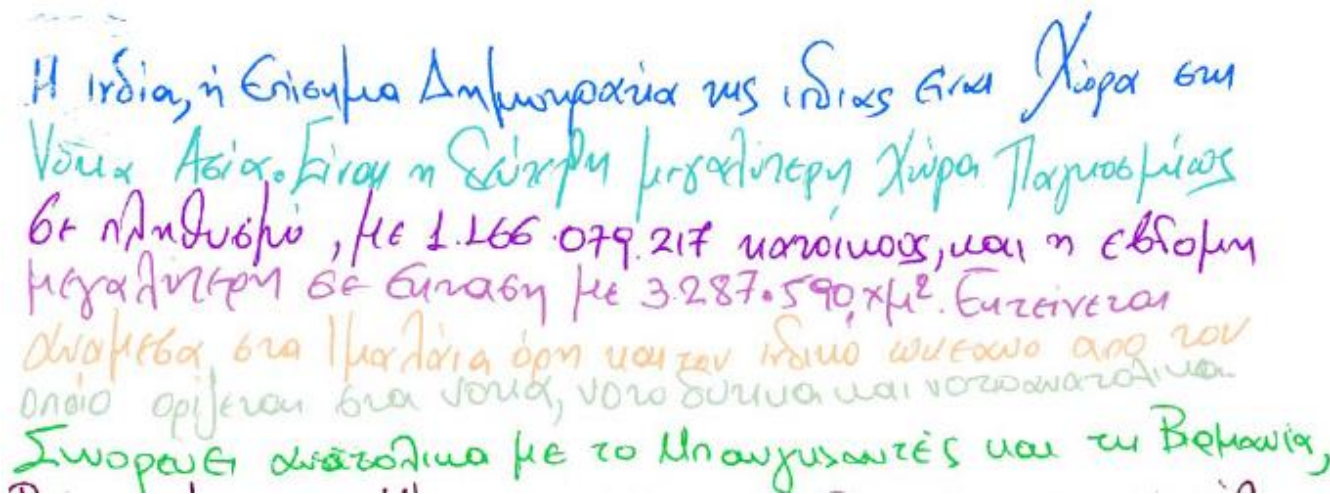
• 版面分析

– Region/text line segmentation

- Top-down: projection, recursive X-Y cut, piecewise projection, whitespace analysis
- Bottom-up: Hough transform, smearing, clustering, Docstrum
- Level set (active contour)
- Seam carving (path finding in energy map)

– Table item extraction

- Rule-based, model based (registration)



Η Ελλάδα, η Ευρώπη Δημοκρατία της Ελλάδας είναι χώρα στη
Νότια Ασία. Είναι η δεύτερη μεγαλύτερη χώρα Παλαιάς
δε ημερομηνία, με 1.166.079.217 κατοίκους, και η εβδόμη
μεγαλύτερη σε έκταση με 3.287.590,4 km². Εμφανίζεται
από το 6ο αιώνα π.Χ. και ήταν η ίδια χώρα από τον
10ο αιώνα π.Χ. και τον 11ο αιώνα π.Χ. και τον 12ο αιώνα π.Χ.
Συνιστάται από το Μοναρχικό και τη Βασιλεία,

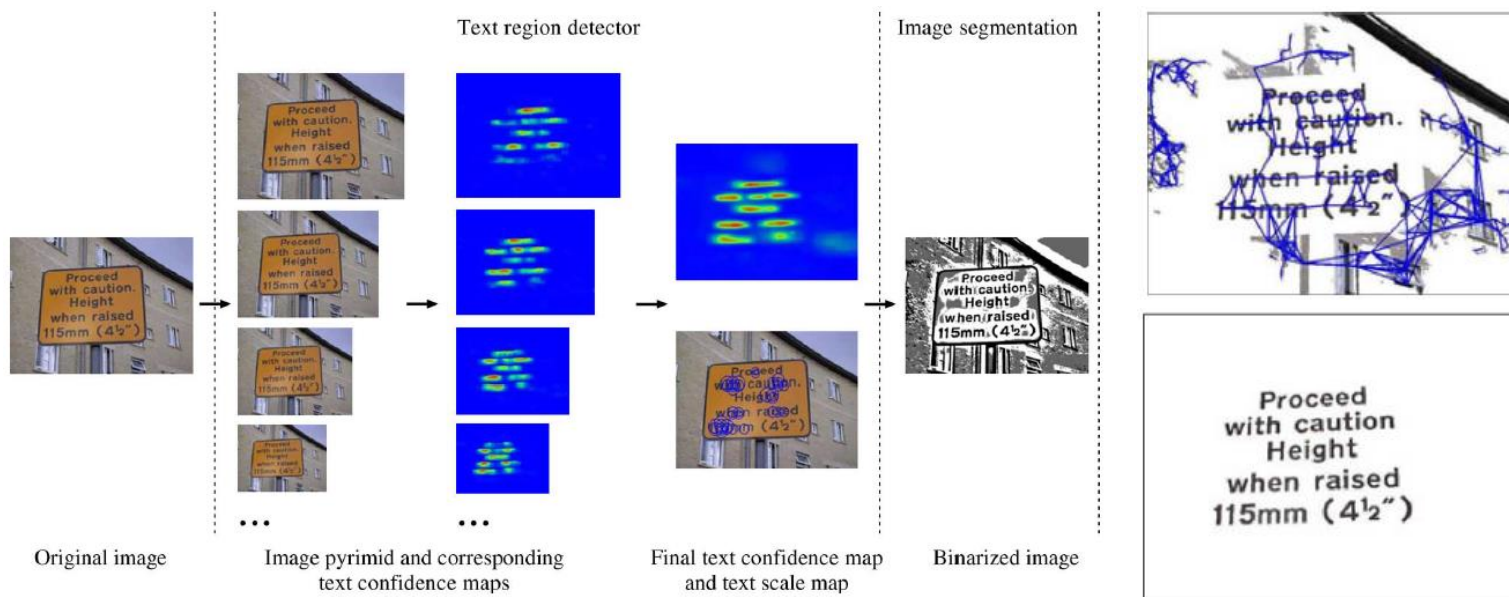
• 版面分析

– Online handwriting

- Text/non-text separation: HMM, MRF, CRF
- Text line segmentation: clustering, dynamic programming

– Text localization

- Region-based: sliding window
- Component-based: SWT, MSER
- Hybrid: component filtering and grouping
- End-to-end: character detection and grouping

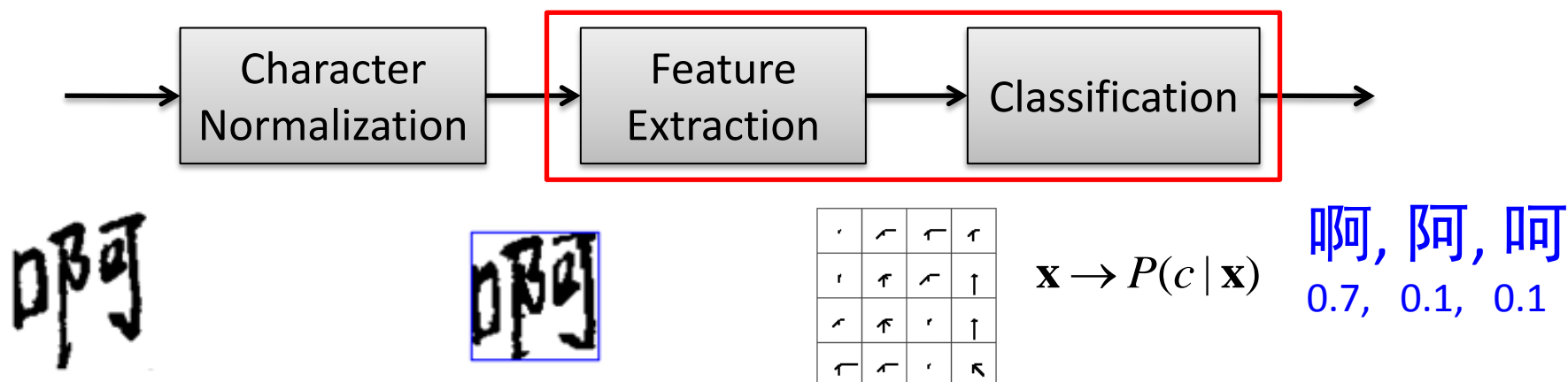


• Text Recognition

– Character recognition

- Normalization: linear, moment-based, nonlinear, pseudo 2D
- Feature extraction: direction histogram, Gabor, structural
- Dimensionality reduction: PCA, FDA, DFE (discriminative)
- Classification: statistical, neural (MLP, RBF, polynomial), SVM
 - Large category set: MQDF, LVQ, hierarchical
- Deep learning

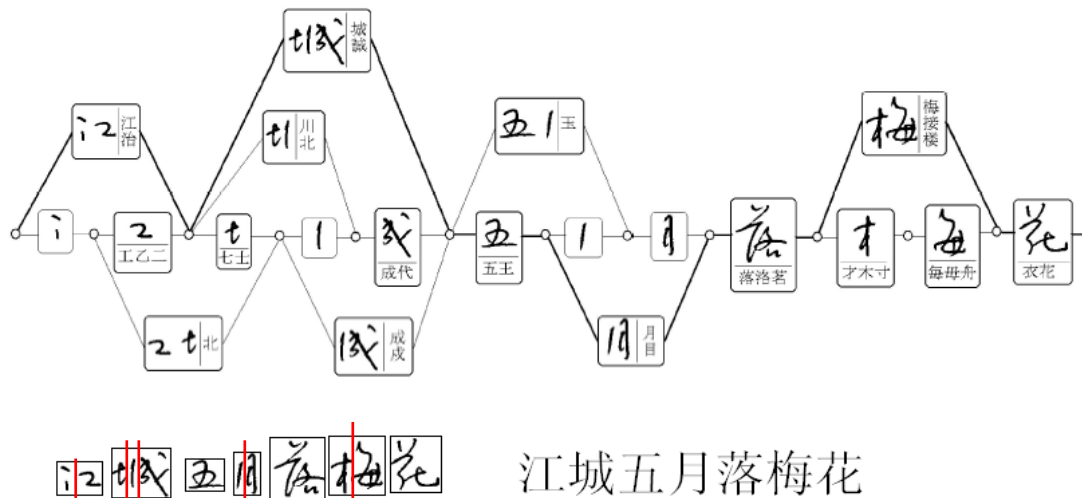
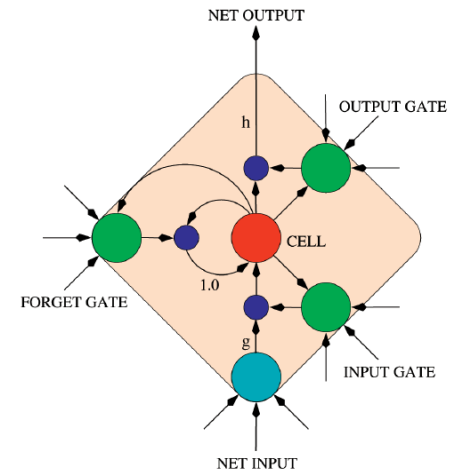
Deep Learning



- Text Recognition

- Word/text line recognition

- Explicit/over segmentation
- Implicit segmentation: sliding window
- HMM
- RNN, BLSTM
(bidirectional long short-term memory)
- Fusion: Bayesian, weighted fusion
- Search: DP, lexicon-free, lexicon-driven
- Language model: n-gram, neural



- Graphics/Symbol
 - Engineering drawings
 - Primitive extraction
 - Graph matching
 - Flowchart
 - Stroke labeling: MRF, CRF
 - Rule-based interpretation
 - Mathematics
 - Symbol segmentation
 - Symbol recognition
 - Graph/grammar/rule-based interpretation
 - Other isolated symbols
 - Statistical, structural

- Style Authentication
 - Writer identification
 - Text dependent/independent
 - Textural features, HOG, LBP
 - Bag of features
 - Nearest neighbor classification
 - Signature verification
 - Statistical, structural
 - Distance metric (intrapersonal, extrapersonal)
 - Font identification
 - Global and local shape features
 - Hierarchical classification
 - Script identification
 - Textural features
 - Key graphemes, n-grams

- Document Understanding
 - Logical ordering: knowledge-based
 - Summarization: word equivalence classes detection, stop words and keywords
 - Categorization: OCRed text classification, keyword spotting based BoW
- Keyword Spotting
 - Query by image/string
 - OCR-based
 - Template matching
 - Shape matching, bag-of-features
 - Model based
 - Character/word classifier, context models
 - Need training

研究现状：性能状况

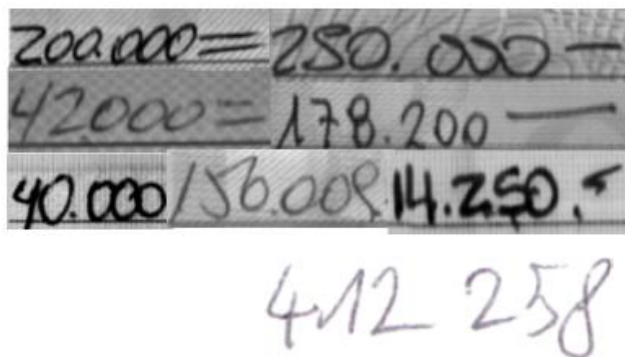
- Handwritten numeral

- Isolated: MNIST

Liu et al. 2003	Error%	#param	Time
MLP	0.60	63.31K	0.44ms
Polynomial	0.58	38.86K	0.76ms
SVC-poly	0.55	913K	5.90ms
SVC-rbf	0.42	1.61M	21.9ms

- Numeral strings

- ICFHR2014 HDSRC



DNN	Error (%)
Simard et al. 2003	0.40
Ciresan et al. 2010 (IDSIA)	0.35
Wu et al. 2014 (Fujitsu)	0.254

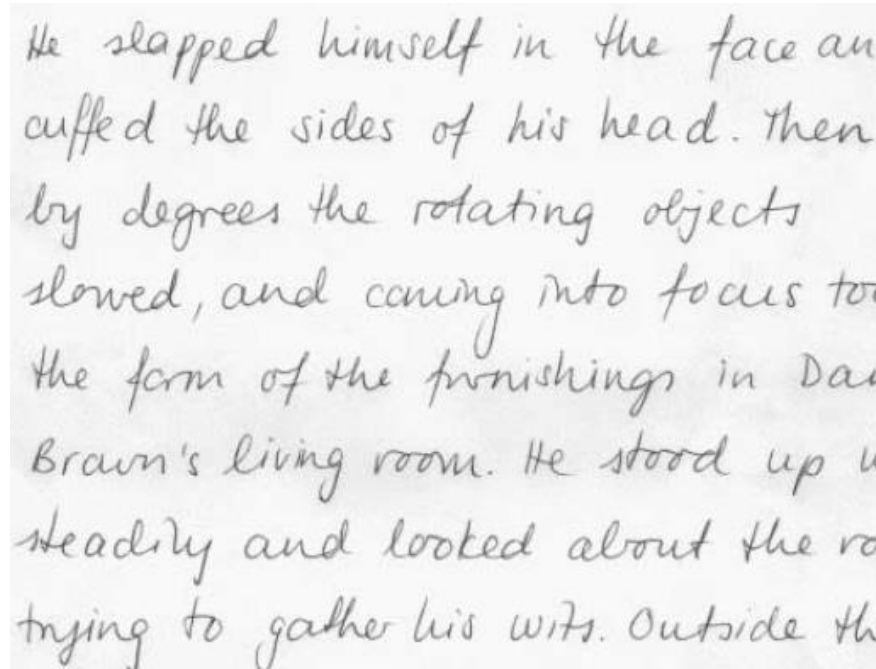
Submission	Guesses	CAR A	CAR B	CVL	Mean
Tébessa I	TOP-1	0.3705	0.2662	0.5930	0.4099
	TOP-2	0.4559	0.3401	0.6575	0.4845
	TOP-3	0.4720	0.3568	0.6690	0.4993
Tébessa II	TOP-1	0.3972	0.2772	0.6123	0.4289
	TOP-2	0.4477	0.3137	0.6527	0.4714
	TOP-3	0.4818	0.3411	0.6824	0.5018
Singapore	TOP-1	0.5230	0.5960	0.5040	0.5410
	TOP-2	0.6180	0.6770	0.6060	0.6337
	TOP-3	0.6540	0.7130	0.6540	0.6737
Pernambuco	TOP-1	0.7830	0.7543	0.5860	0.7078
	TOP-2	0.8916	0.8746	0.6850	0.8171
	TOP-3	0.9199	0.9009	0.7234	0.8481
Beijing	TOP-1	0.8073	0.7013	0.8529	0.7872
	TOP-2	0.8634	0.7638	0.9128	0.8467
	TOP-3	0.8697	0.7779	0.9189	0.8555
Shanghai	TOP-1	0.4950	0.2809	0.4893	0.4217
	TOP-2	0.5378	0.3120	0.5400	0.4633
	TOP-3	0.4950	0.2809	0.4893	0.4849

- Handwritten English Texts

- IAM (Univ. Bern) database

- 13,040 lines, containing 86,272 instances of 11,050 distinct words
 - 6,161 lines for training, 2,781 lines for testing

	Method	WER %	CER %
Graves et al. 2009	RNN BLSTM	25.9	18.2
Espana et al. 2011	HMM/ANN	21.14	9.1
Ney et al. 2013	Tandem HMM	13.3	5.1
Ney et al. 2014	RNN-HMM	12.2	4.7



He slapped himself in the face and
cuffed the sides of his head. Then
by degrees the rotating objects
slowed, and coming into focus took
the form of the furnishings in David
Brown's living room. He stood up un-
steadily and looked about the room
trying to gather his wits. Outside the

- Handwritten Arabic Texts
 - NIST OpenHaRT'13
 - Three tasks: DIR, DIT, DTT
 - Given text line segmentation
 - Constrained condition: given training set
 - Unconstrained condition

DIR results (1-WER)

Team ID	Method	Constr	Unconstr
A2iA (Fra)	RNN	0.7990	0.8164
RWTH (Ger)	HMM	0.7624	0.8390
CITLAB (Ger)	RNN	0.7373	
UPV (Spa)	HMM	0.7068	
UOB-TPT (Fra)	RNN-HMM	0.5207	
LITIS (Fra)	HMM	0.2241	

انفجرت سيارة مفخخة أمس بالقرب من
السيارة الدبلوماسية وأحد مدخل المنطقة
التي فيها الشرطة التحصين في وسط
خوار خمسة آخرين بجرح، فيما أعلنت
مصادر الخفاضة معدلات القتلى في صفوف
الفراتيين والعسكريين خلال الشراء
وقال مصدر في وزارة الداخلية العراقية إنه
الانفجار وقع في ساحة لوقوف الطائرات قرب
مدخل يؤدي إلى المنطقة الخضراء يستخدم
الحامول في وزارة الدفاع القريبة من
المكان، وقرب مقر السفارة العراقية -
ومن ضمنه أخرى كوزارة الخارجية الفرنسية
وإدارة كوشنير في بغداد، أمس، إنه
الانفجار الرئيسي تشدد حذرا من التحسن من
والسيارة التي تسير في الحالة «الخطيرة» ما في

القدس 16 - 6 - 2008 (أ ف ب) -
صارت لجنة تخطيط مبنى إسرائيل الجديد
على خطط لبناء أربعين ألف وحدة سكنية
حلال العهد المقبل في القدس، بعضها في
أحياء استيطانية في القدس الشرقية
المحتلة، على ما أفادت بلدية القدس
اللاتينية.
وسمى بناء قسري هذه المساكن في أحياء
في القدس الغربية، كما تنص الخطط على
قيام مقاولين مع القطاع الخاص لبناء آلاف
المساكن لسكان القدس الشرقية الفلسطينيين
اطقم خدام بنحو مئتي ألف نسمة.
وصارت لجنة التخطيط الهدي في القدس
التابعة لوزارة الداخلية على الخطم بجرمها
هادقت عليها البلدية.
ورفضت الطرحة ببناء البلدية ونسب غوتسمان
مرا على أسئلة قرأه بركن تحديد عدد
المساكن التي ستشيد في القدس
الشرقية، موضعا أنه البلدية «لا تفكر

• Handwritten Chinese Characters

– CASIA OLHWDB/HWDB

躲朵 躲朵 躲朵 躲朵 躲朵 躲朵 躲朵 躲朵 躲朵
 佛 额 沁 娥 恶 厄 扼 遏 鄂 饿
 恩 而 儿 耳 尔 洱 二 贰 发
 罚 筏 伐 乏 阀 法 珐 藩 帆 番
 翻 樊 矾 钒 繁 凡 烦 反 返 范
 贩 犯 饭 泛 坊 芳 方 防 房 防
 妨 仿 访 纺 放 菲 非 啡 飞 肥
 匪 啡 味 肺 废 沸 费 芬 酚 吩

躲朵 躲朵 躲朵 躲朵 躲朵 躲朵 躲朵 躲朵 躲朵
 佛 额 沁 娥 恶 厄 扼 遏 鄂 饿
 恩 而 儿 耳 尔 洱 二 贰 发
 罚 筏 伐 乏 阀 法 珐 藩 帆 番
 翻 樊 矾 钒 繁 凡 烦 反 返 范
 贩 犯 饭 泛 坊 芳 方 防 房 防
 妨 仿 访 纺 放 菲 非 啡 飞 肥
 匪 啡 味 肺 废 沸 费 芬 酚 吩

– ICDAR 2013 competition

Table 4. Results of online character recognition (%).

System	CR (1)	CR (10)	Ave time	Dic size
UWarwick	97.39	99.88	355ms	37.8M
VO-3	96.87	99.67	15.3ms	87.6M*
VO-2	96.72	99.61	4.1ms	36M*
VO-1	96.33	99.61	1.6ms	10M*
HIT	95.18	99.39	2.3ms	120M
USTC-2	94.59	99.14	3.8ms	5.25M
USTC-1	94.25	99.06	2.0ms	3.19M
TUAT	93.85	99.24	5.3ms	96.2M
Faybee	92.97	98.87	0.5ms	4.48M
Ref [1]	95.31			
Human	95.19			

Offline character recognition

	System	Error (%)	Speed (ms)
ICDAR2013 Competition	Fujitsu, CNN	94.77	55 (GPU)
	IDSIAAnn (8)	94.42	315 (CPU)
	IDSIAAnn (4)	94.24	197 (CPU)
	HIT, traditional	92.62	4.6 (CPU)
IDSIA Tech Rep 05-13 (2013)	CNN	94.47	3.03 (GPU)
	Multi-CNN (8)	95.78	22.04 (GPU)
Fujitsu (ICFHR2014)	ATR-CNN	95.04	
	CNN voting	96.06	

- Handwritten Chinese Texts

- ICDAR2013 competition: given text line segmentation

中医认为,痤疮患者大多数有内热,应多食一些 瘦猪肉,猪肉、兔肉、鸭肉、鱼翅鱼、蘑菇、银耳、黑木耳、芹菜、菠菜、苋菜、莴笋、苦瓜、丝瓜、冬瓜、黄瓜、西瓜、西红柿、绿豆、绿豆芽、黄豆芽、豆腐、莲藕、梨、桑椹、柚子、山楂、苹果等,这些食物有起清凉去热、生津润燥的作用。中医认为,痤疮患者主要是过食肥甘厚味,导致肺胃湿热熏蒸,面部肌肤所引起。因此,凡含油脂丰富的食品,如肥肉、动物脑、蛋黄、芝麻、花生等,都应少吃。中医认为,辛辣湿热食物,如烟、酒、浓茶、咖啡、辣椒、大蒜、韭菜、狗肉、雀肉、虾等,会使痤疮加重或复发,应忌食。

Table 5. Results of offline text recognition (%).

	CR	AR	Ave time	Dic size
HIT-2	88.76	86.73	1.2s	309M
HIT-1	86.15	83.58	0.64s	111M
THU	82.92	79.81	0.85s	102M
SCUEC	42.05	35.14	0.15s	442M
Ref[6]	90.22	89.28		

Table 6. Results of online text recognition (%).

	CR	AR	Ave time	Dic size
VO-3	95.03	94.49	1.72s	56M*
VO-2	94.94	94.37	1.23s	37.9M*
VO-1	93.11	92.57	0.72s	20.8M*
TUAT	88.49	87.66	1.42s	246M
USTC	82.20	81.57	0.25s	29.3M*
Ref [29]	94.62	94.06		

*Size of executive file embedding dictionary.

- Scene Texts
 - ICDAR2013 Robust Reading Competition

Text localization results

Method Name	Recall (%)	Precision (%)	F-score
USTB_TexStar	66.45	88.47	75.89
Text Spotter [20], [21], [22]	64.84	87.51	74.49
CASIA_NLPR [23], [24]	68.24	78.89	73.18
Text_Detector_CASIA [25], [26]	62.85	84.70	72.16
I2R_NUS_FAR	69.00	75.08	71.91
I2R_NUS	66.17	72.54	69.21
TH-TextLoc	65.19	69.96	67.49
Text Detection [15], [16]	53.42	74.15	62.10
Huo et al. 2014	85.72	87.03	86.37

Word recognition results

Method	Total Edit Distance	Correctly Recognised Words (%)
PhotoOCR (Google)	122.7	82.83
PicRead [27]	332.4	57.99
NESP [19]	360.1	64.20
PLT [18]	392.1	62.37
MAPS [17]	421.8	62.74
Feild's Method	422.1	47.95
PIONEER [28], [29]	479.8	53.70
Baseline	539.0	45.30
TextSpotter [20], [21], [22]	606.3	26.85



Results on Street View Text dataset (ICCV2013)

Algorithm	Word Recognition Rate (%)
PhotoOCR	90.39
Goel et al. [9]	77.28
Mishra et al. [15]	73.26
Novikova et al. [20]	72.9
Wang et al. [26]	70.0
Baseline (ABBYY) [9]	35.0

国内现状

- 1990-1999是黄金时期
 - 863评测（3次以上）
 - 研究组：中科院自动化所（刘迎建），清华大学（丁晓青、夏莹），北京大学（顾晓凤），中科院计算所（张永慧），上海交大（施鹏飞），哈工大（唐降龙），武汉工业大学（胡家忠）
- 2000年前后跌入低谷
- 2005年以后逐渐恢复
 - 跨国公司（富士通、微软）重视，基金支持，智能手机带来机会
 - 中科院自动化所，清华大学，华南理工大学（金连文），华东师范大学（吕岳）
 - Camera-based: 百度，北京科技大学（殷绪成），华中科技大学（白翔），微软（霍强），南京大学（路通），中科院自动化所（孟高峰），西安交大（宋永红）
 - 北大计算机所：字体，PDF文档，数学公式
- 国际影响
 - ICPR: 2010年以来中国学者论文数居第一
 - ICDAR: 2011中国与美国录用论文数相当，此外与第一名差距较大

研究趋势

- Major Challenges
 - Unconstrained handwriting
 - Isolated characters
 - Continuous handwriting
 - Online handwritten notes: mixed texts with graphics
 - Camera-based OCR
 - Scene text: illumination, perspective, cluttered background
 - Camera-based documents: illumination, perspective
 - Historical documents
 - Degraded image
 - Large character set, many variants of same class
 - Few labeled samples
 - Multilingual documents
 - Mixed languages, unknown a priori

研究方向

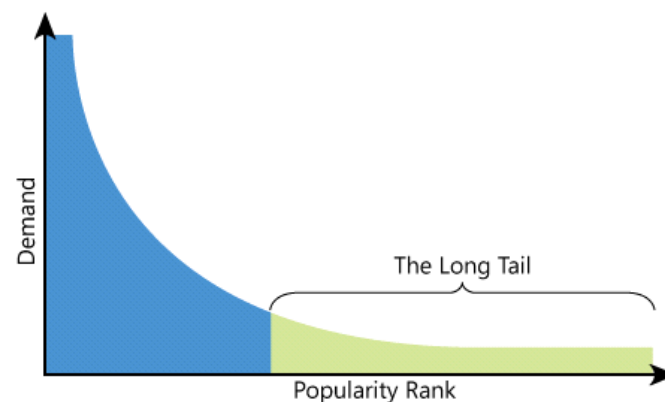
- Fundamental Theory & Methodology
 - Learning for large category/scale data, structured model, weakly labeled data (deep learning)
 - Fusion of multi-level/source contexts, global optimization
 - Online learning
 - Cognitive mechanisms: from visual cues to high-level knowledge
- Document Image Processing
 - Image capturing paradigm/device
 - Camera-captured image rectification
 - Binarization of degraded image
 - Complex layout analysis

- Isolated Character Recognition
 - Important for string recognition, higher accuracy always desired, non-char rejection
 - Feature extraction/learning
 - Classifier learning and adaptation
 - Confidence/reliability
- Text Line Recognition
 - Sequence classification model and learning
 - Touching character segmentation
 - Contexts modeling and fusion
 - Language modeling and adaptation
 - Global optimization of integrated string model
 - Multi-lingual documents, especially mixed languages

- Scene Text Recognition
 - Vision computing inspired feature and detector
 - Text extraction (separation from background)
 - Scene text string recognition
 - Principled fusion of bottom-up and top-down cues
- Application Oriented
 - Fast implementation, light implementation
 - Modeling of interactive transcription
 - Data generation and annotation
 - Document retrieval and information extraction
 - Document authentication, writer identification
 - Online handwritten notes
 - Application to human interface, robot, archeology, education, impaired person assistance, etc

展望

- 领域发展的动力
 - 应用需求
 - 挑战性
- 文档图像：存在就是需求
 - 文档何时会消亡？永远不会
 - 技术何时能完全解决？至少20年内不会
- 需求：Long Tail
 - 敏感信息多存在于图像中
 - 李彦宏：KDD2012 Keynote
- Document Image is Big Data
 - 3V: volume, velocity, variety



探讨

- 深度学习是否会代替一切
 - 深度学习只是部分解决问题，而且有代价
 - 字符识别、字符串识别上，传统方法仍然可能与深度学习竞争（提高特征维数、增强训练样本）
 - 传统方法和深度学习都需要进一步发展
- 学术界如何生存
 - 如何找学术问题
 - 技术问题后面一定有学术问题，能发表高档次论文
 - 与工业界互动
 - 学术界（前瞻、深入）与工业界（应用）定位不同
 - 发表论文就不能保护技术，学术与应用难以兼顾
 - 合作的前提是互补，相互有利用价值

探讨

- 如何吸引年轻人进入DA领域
 - 就业情况：不差
 - DA研究基础也能在相关领域（如DM, CV）就业
 - 发表论文：有很多机会
 - 入门不难，但深入很难
 - 研究条件：数据平台、实验平台（已有技术模块）
 - 开源代码：DA领域严重不足
 - 如何开源同时又保护自己的技术？
 - 不能产品化的代码：Matlab, DLL, 不是领先但是常用的算法

Acknowledgements

- Statistics of Past Conferences
 - Ching Suen, Hiromichi Fujisawa, Masaki Nakagawa, Andreas Dengel, Daniel Lopresti

**Thank you for
attention!**