# NAVIGATING SCIENTIFIC LITERATURE
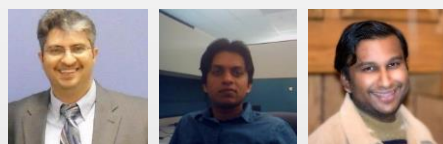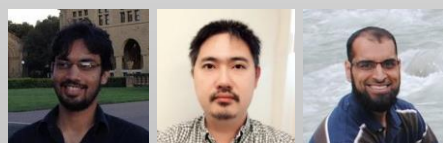## A HOLISTIC PERSPECTIVE

Venu Govindaraju

## PATTERN RECOGNITION

Towards a Globally Optimal Approach for Learning Deep Unsupervised Models

Organizing Multiple Experts for Efficient Pattern Recognition

Active Pattern Recognition Using Genetic Programming

A Complexity Framework for Combination of Classifiers in Verification and Identification Systems

Image Processing using Ontology Concepts for Image Segmentation

Language Motivated Approaches for Human Action Recognition and Spotting

## DOCUMENT ANALYSIS

Intrusion Detection using Spatial Information and Behavioral Biometrics

Integrating Minutiae Based Fingerprint Matching with Local Correlation Methods

Integrating Facial Expressions and Skin Texture in Dace Recognition

Stochastic Modeling of High-level Structures in Handwritten Word Recognition

Statistical Techniques for Efficient Indexing and Retrieval of Document Images

Probabilistic Random Field based Text Identification

Enhancing Cyber Security through the use of Synthetic Handwritten CAPTCHAs

Language Models and Automatic Topic Categorization for Information Retrieval in Handwritten Documents

Methods for Biomedical Image Content Extraction Toward Improved Multimodal Retrieval of Biomedical Articles

A Novel Multi-sample Fusion Methodology for Improving Biometric Recognition

Enhancement and Retrieval of Low Quality Handwritten Documents

A Stochastic Framework for Font Independent Devanagari OCR

A Semi Supervised Framework for Handwritten Document Analysis

Bayesian Background Models for Retrieval of Handwritten Documents

Accents in Handwriting: A Hierarchical Bayesian Approach to Handwriting Analysis

Hierarchical and Dynamic-Relational Models for Handwriting Recognition

## BIOMETRICS

Multilingual Word Spotting in Offline Handwritten Documents

A Framework for Fingerprint Enhancement and Feature Detection

Minutia-Based Partial Fingerprint Recognition

Sequential Pattern Classification without Explicit Feature Extraction

Automatic Recognition of Handwritten Medical Forms for Search Engines

Exploiting the Gap between Human and Machine Abilities in Handwriting Recognition for Web Security Applications

Face Modeling and Biometric Anti-spoofing using Probability Distribution Transfer Learning

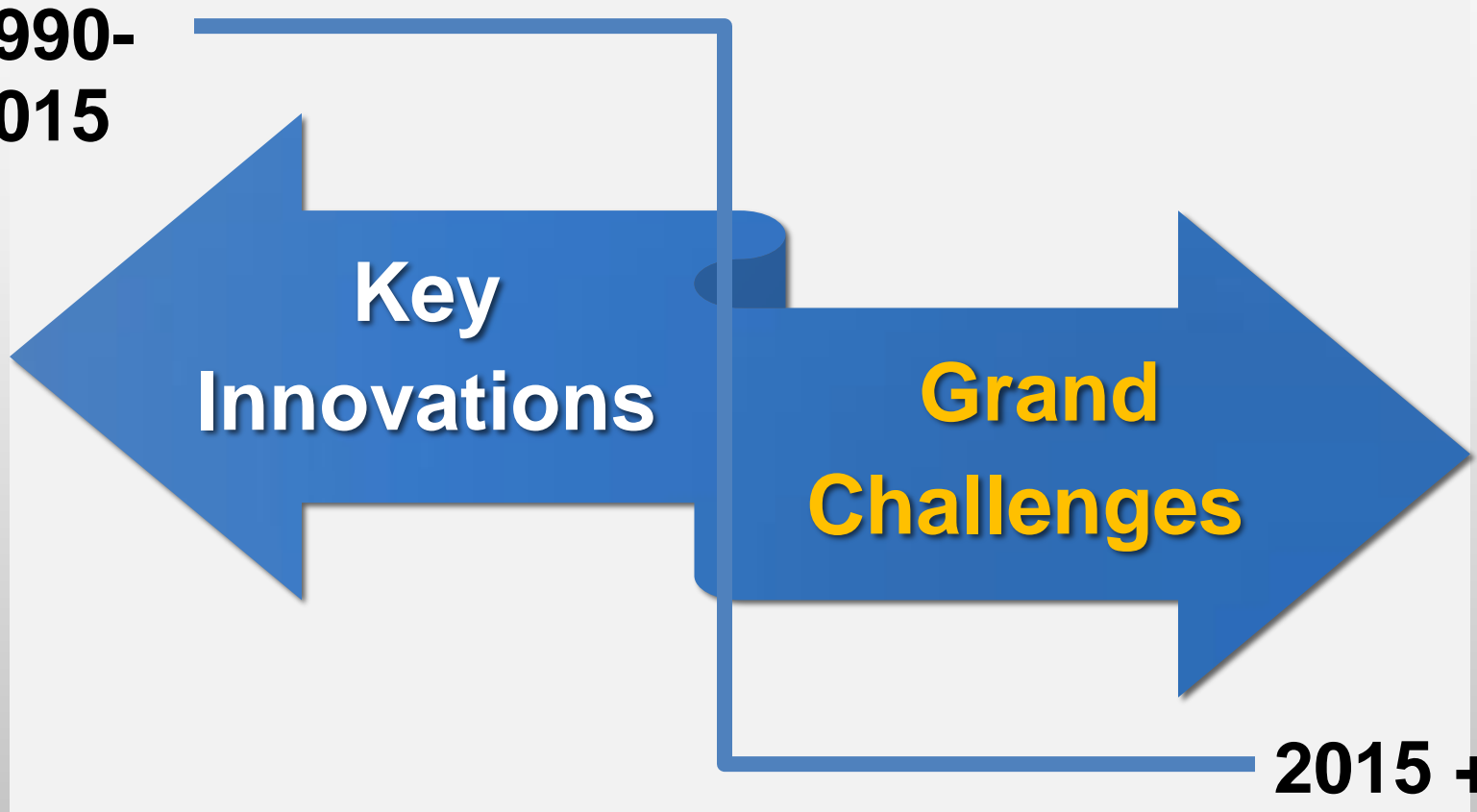A Framework for Efficient Fingerprint Identification using a Minutiae Tree

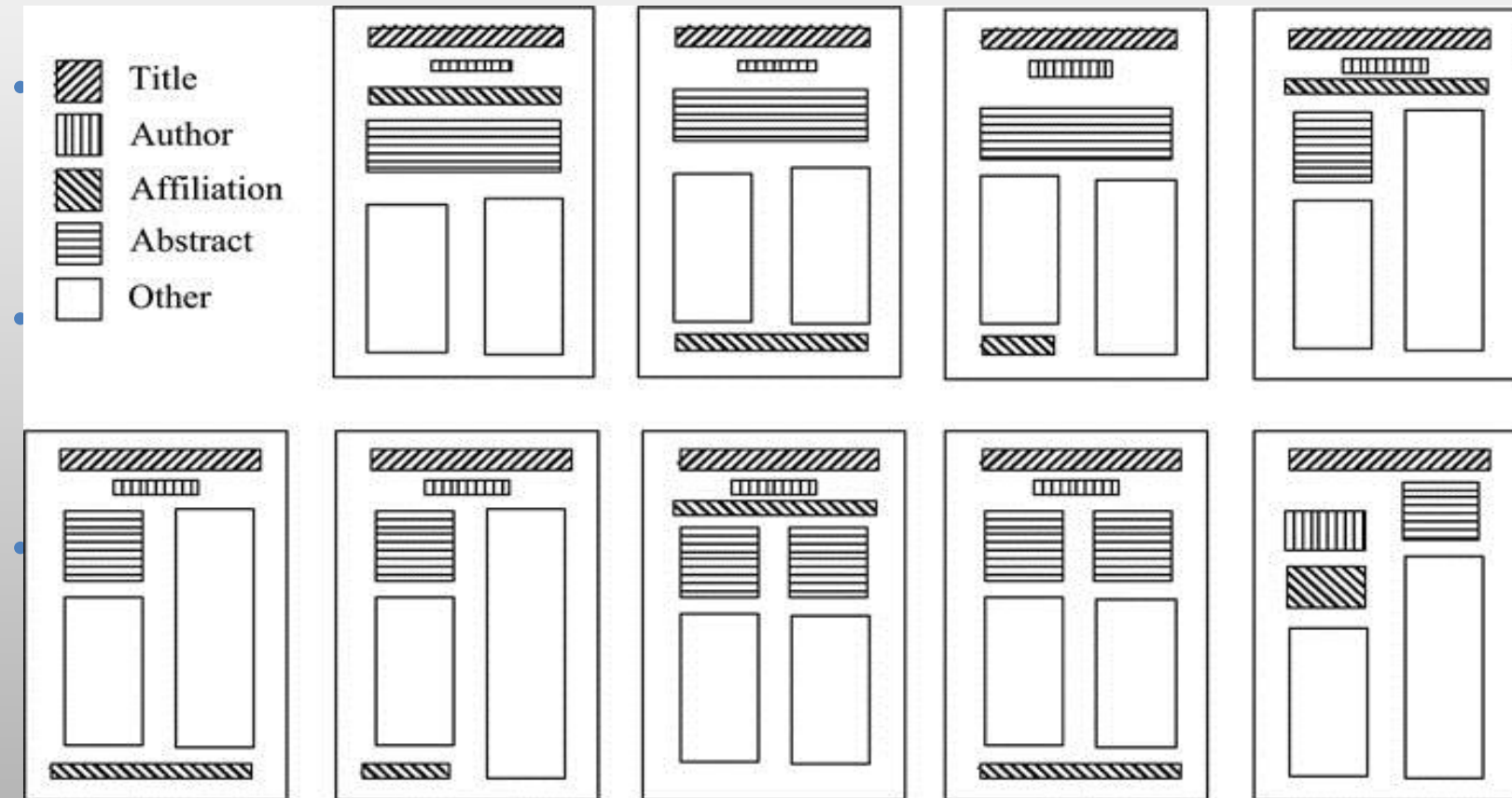33

1990-2015

Key Innovations

Grand Challenges

2015 +

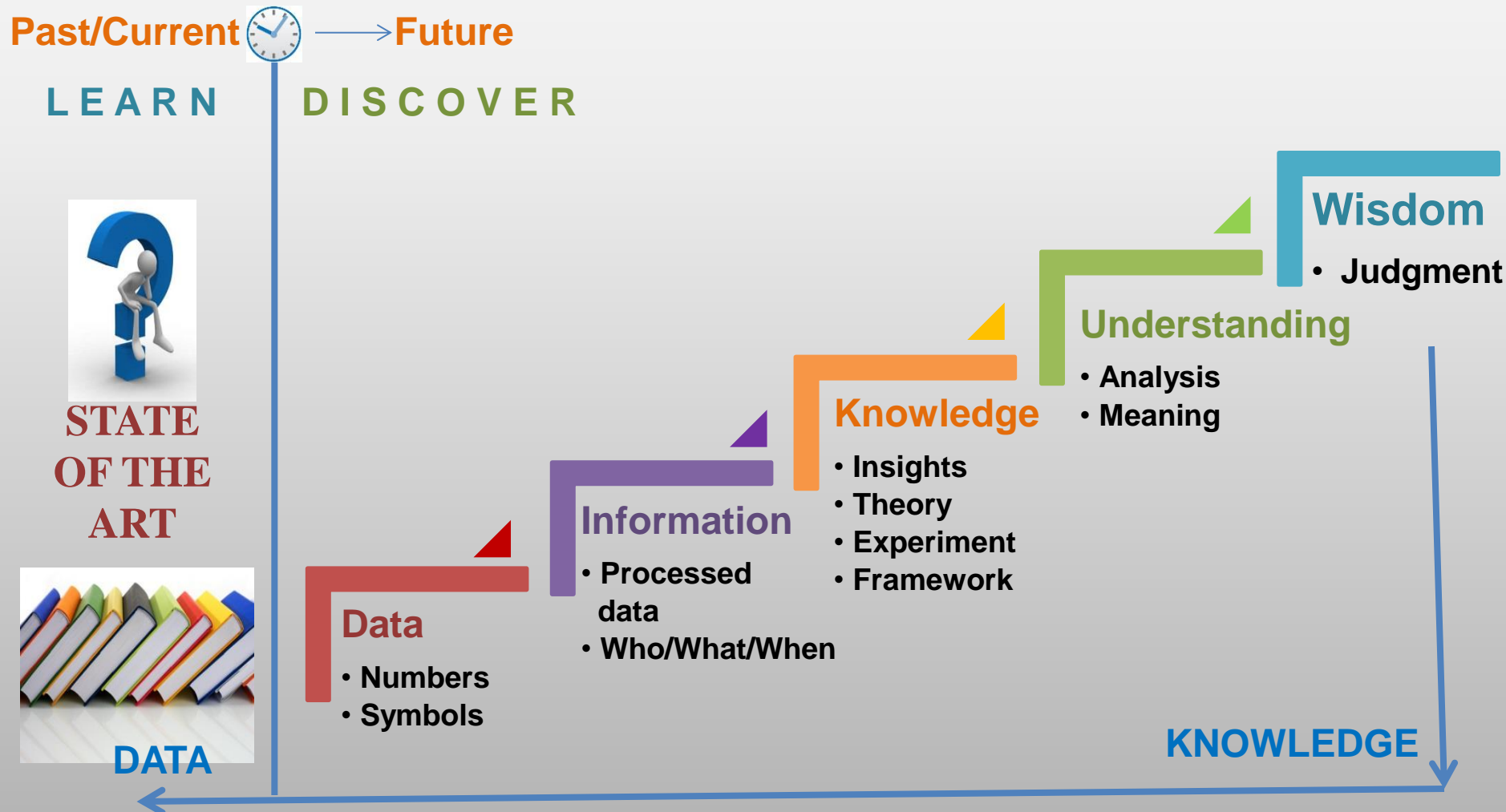# Old Order - DIA

UW English Document Image Database
(Phillips, Technical report, 1996, citations: 29)

# Scientific Process
### nanos gigantum humeris insidentes

1676 letter of Isaac Newton: " If I have seen  further it is by standing on the shoulders of giants."

**Past/Current** → **Future**

**L E A R N**     **D I S C O V E R**

**STATE OF THE ART**

**Wisdom**
- **Judgment**

**Understanding**
- **Analysis**
- **Meaning**

**Knowledge**
- **Insights**
- **Theory**
- **Experiment**
- **Framework**

**Information**
- **Processed data**
- **Who/What/When**

**Data**
- **Numbers**
- **Symbols**

**DATA**

**KNOWLEDGE**

# When knowledge becomes data



Variety

Velocity

**Big Data**

Volume

Veracity

# Scientific Literature

- 2009 estimate: 50 million articles;  28 thousand journals
- 1.8M articles added every year.

| | | |
|---|---|---|
| PubMed — A service of the U.S. National Library of Medicine and the National Institutes of Health www.pubmed.gov | 23 million articles (Just biomedical literature) | **Volume** |
| Scopus — The largest abstract and citation database of peer-reviewed literature. | 45 million articles | **Variety** |
| THOMSON REUTERS WEB OF KNOWLEDGE | 40 million articles | |
| Google scholar | Unknown (peer reviewed only) | **Veracity** |

BIG DATA!

Roughly, papers double every 10-15 years !                    **Velocity**

[Meadows, 1998, p.16]

# Big Data Side Effects
## Challenges for ICDAR Community

4Vs

| | |
|---|---|
| Volume | References |
| Velocity | Reinvention |
| Variety | Replicability |
| Veracity | Reputation |

4RS

# References
## Volume Challenge

The Royal Society (1662)
First journal : *Journal of Philosophical Transactions (1665)*
*Le Journal des Sçavans (1665)*

Scientists believe they are only reading 40% of the relevant literature.
Faraday reported the same problem already in 1826 !!

[Meadows, 1998], page 211, and Faraday is quoted on page 19

"50% of papers are never read by anyone other than their authors, referees and journal editors."….

"90% of papers are never cited …"

[Smithsonian.com, 2007 study]

# **R**einvention
## **V**elocity Challenge

``in some disciplines it is occasionally easier to repeat an experiment than it is to determine that the experiment has already been done.'' [Garvey, 1979, p.8].

# **R**eplicability
## **V**eracity Challenge

- ***Nullius in verba***
  "On the word of no one" or "Take nobody's word for it"

SCIENCE is in crisis, just when we need it most. Two years ago, C. Glenn Begley and Lee M. Ellis reported in Nature that they were able to replicate only six out of 53 "landmark" cancer studies. Scientists now worry that many published scientific results are simply not true.

**NY Times 2014**

# **R**eputation
## **V**eracity Challenge

- How Many Scientists Does It Take to Write a Paper?

Scientific journals see a spike in number of contributors; 24 pages of alphabetized co-authors.

**The Wall Street Journal, August 10, 2015**

# Challenges

# GRAND   CHALLENGE

## BIG DATA

**VOLUME**

References

Cognitive burden

**VELOCITY**

Reinvention

Reinventing wheel

**VARIETY**

**VERACITY**

Replicability

Authenticity

# Addressing the Cognitive Burden
**Volume**

# Addressing Reinventing the wheel ?
## Velocity

- **Least square with linear constraints:** one type of quadratic program in mathematics

$$\text{minimize} \quad \|Ax - b\|_2^2$$
$$\text{subject to} \quad l_i \leq x_i \leq u_i, i = 1, \ldots, n$$

- **Isotonic regression:** in statistics

$$\text{minimize} \quad \sum_{i=1}^{n} w_i(x_i - a_i)^2$$
$$\text{subject to} \quad x_i \geq x_j \, for \, (i,j) \in E$$



**Trapezoid rule: calculus 17th century**



Figure 1—*Total area under the curve is the sum of individual areas of triangles a, c, e, and g and rectangles b, d, f, and h.*

**Tai's Model, 1994, 254 citations**[17]

# Addressing **Replicability**
## Velocity



IBM Journal 1982

- **Dataset** – UNLV/ISRI
  - 64 pages, 6796 blocks

- **Heuristics parameters**
  - Vertical: 15 pixels
  - Horizontal: 50/30 pixels

- **Classes:**
  - Text, Table, Caption, Figs
- **Classifier:**
  - Support Vector Machines

- **Accuracy:**
  - 91.73 % at block level

# Addressing Authenticity
## Veracity

### Datasets

- Public

- Benchmark

- Published

### Reputation

- Authors

- Lab

- Journal

### Experiments

- Comparative results

- CODE available !!

### Citations

- Only Counts ?

# Veracity
## All citations are not equal

▪ **Which citation is more trustworthy?**

object classification [3]. However, the same level of success has not been obtained for generative tasks, despite numerous efforts [13, 24, 28].

Table 2. Results on MNIST dataset.

| Method | Paper | Error rate[%] |
|---|---|---|
| CNN | [32] | 0.40 |
| CNN | [26] | 0.39 |
| MLP | [5] | 0.35 |
| CNN committee | [6] | 0.27 |
| MCDNN | this | **0.23** |

Sentiment analysis:
Targeted NLP

We are unable to replicate the results from paper [14]

area of speech recognition, with breakthrough results (Dahl et al., 2010; Deng et al., 2010; Seide et al., 2011a; Mohamed et al., 2012; Dahl et al., 2012; Hinton et al., 2012) obtained by several academics as well as researchers at industrial labs

# Veracity
## Dataset linkages

MNIST:   60k training, 10k testing images
"Gradient-based learning applied to document recognition", Lecun et al 1998  (Citations: 3547)

Ciresan et al. 2012

descent with an annealed learning rate. During training, images are continually translated, scaled and rotated (even elastically distorted in case of characters), whereas only the original images are used for validation. Training ends once the validation error is zero or when the learning rate reaches its predetermined minimum. Initial weights are drawn from a uniform random distribution in the range [−0.05, 0.05].
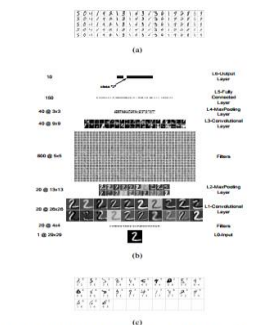
Table 2. Results on MNIST dataset.

| Method | Paper | Error rate[%] |
| --- | --- | --- |
| CNN | [32] | 0.40 |
| CNN | [26] | 0.39 |
| MLP | [5] | 0.35 |
| CNN committee | [6] | 0.27 |
| MCDNN | this | **0.23** |

Training on automatically augmented dataset:
"During training the digits are randomly distorted …
The MCDNN has a very low 0.23% error rate"

# OUR  NEXT FRONTIER

Tables, Graphs

Equations

Targeted NLP

Keyword spotting

Multimedia

# Tables Analysis

**4Vs**



- Extract data
- Compare data
- Headers
- Merged columns and rows
- Color
- Caption?
- Reference in text?

- Such tables in other articles?

# Graphs Analysis
**4Vs**



Figure 8: ROC OF HUMAN AND COMPUTER PERFORMANCE ON MATCHING FACES ACROSS ILLUMINATION CHANGES. ROCs FOR ALGORITHMS IN FIGURE 7 ARE PLOTTED. THE ROC PLOTS FAR AGAINST FRR. PERFECT PERFORMANCE WOULD BE THE LOWER LEFT HAND CORNER (FAR=FRR=0).

- Type of graph? x-y plot, bar graph etc.
- Labels on the axes?
- Number of curves?
- Color?
- Curves intersect?
- Actual data points?
- Legend?
- Figure number?
- Caption?
- Reference in text?
- Such graphs in other articles?

# Equations Analysis
## Velocity

$$y = X\beta + \varepsilon$$

| Domain awareness | |
|---|---|
| 🟥 | Matrix representation |
| ⬛ | Operators |
| 🟦 | Symbols |

| Document awareness |
|---|
| Dependent variables independent variables regression coefficients, error |

$$h_i = \beta_1 t_i + \beta_2 t_i^2 + \varepsilon_i$$

$$y_i = \beta_1 x_{i1} + \cdots + \beta_p x_{ip} + \varepsilon_i$$

$$\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} x_{11} & \cdots & x_{1p} \\ x_{21} & \cdots & x_{2p} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{np} \end{pmatrix} \begin{pmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

# Query Interface

Face Recognition FAR (0.0-0.2) vs FRR

☑ CVPR  ☐ Science  ☐ Nature

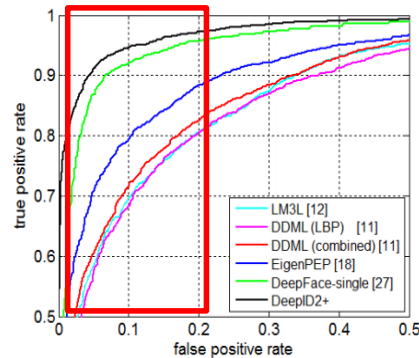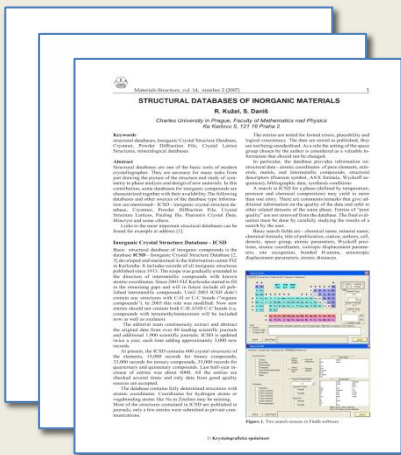☑ Advanced Search Options

☑ X-Y Plots  ☐ Tables  ☑ Figures



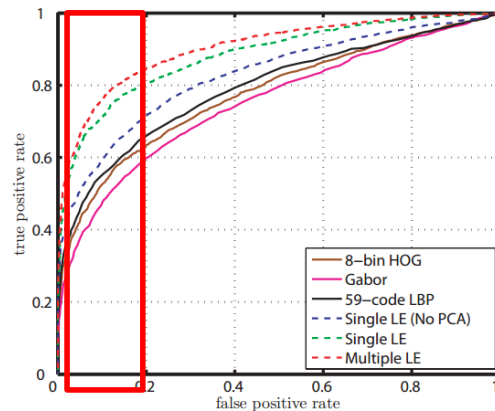Figure 5: ROC of face verification on YouTube Faces. Best viewed in color.

Original Paper



Figure 7. ROC curve comparison between our LE descriptors and existing descriptors.
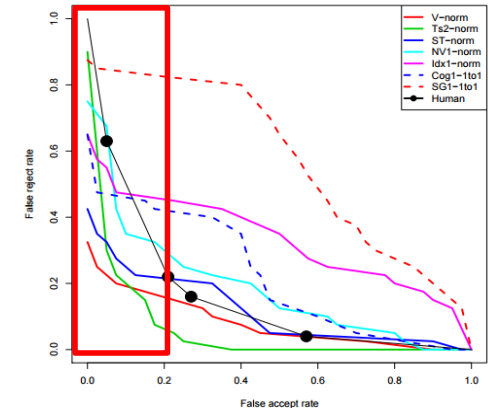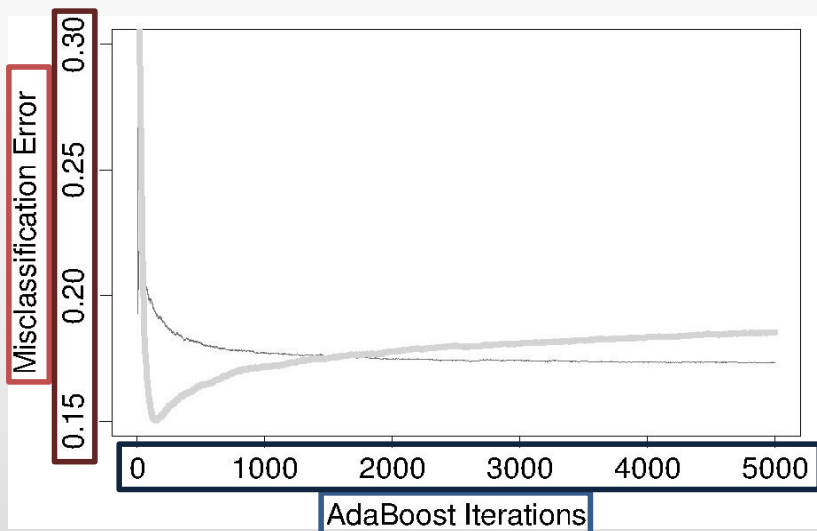
Original Paper



Figure 8: ROC OF HUMAN AND COMPUTER PERFORMANCE ON MATCHING FACES ACROSS ILLUMINATION CHANGES. ROCS FOR ALGORITHMS IN FIGURE 7 ARE PLOTTED. THE ROC PLOTS FAR AGAINST FRR. PERFECT PERFORMANCE WOULD BE THE LOWER LEFT HAND CORNER (FAR=FRR=0).
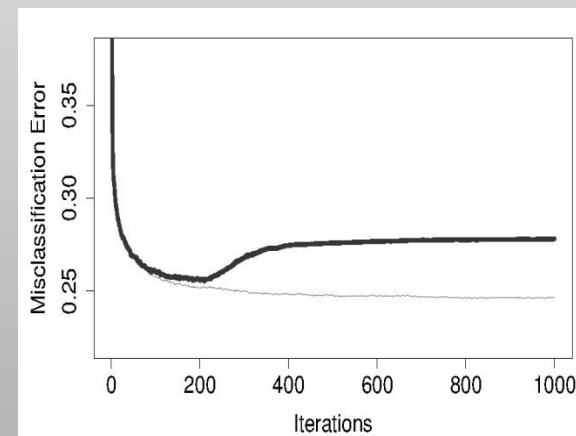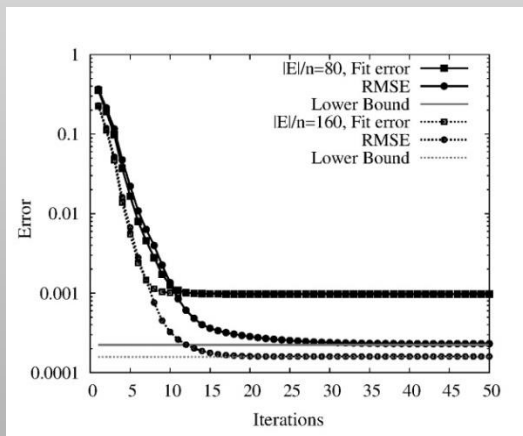
Original Paper

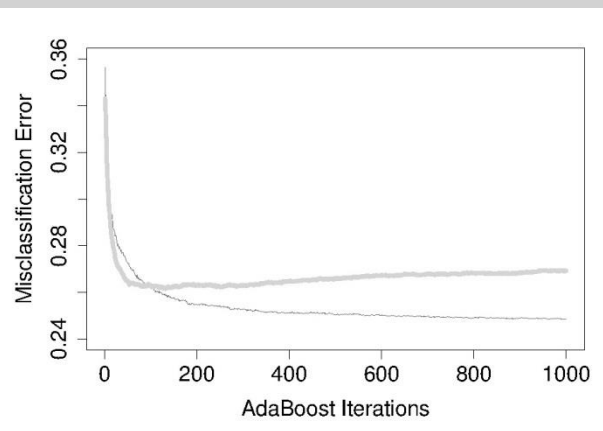## Query:



## Retrieve similar graphs

Line graph

- X Axis
  - Label: AdaBoost iterations
  - Range: 0-5000
  - Unit: -

- Y Axis
  - Label: Misclassification Error
  - Range: 0.15-0.30
  - Unit: -
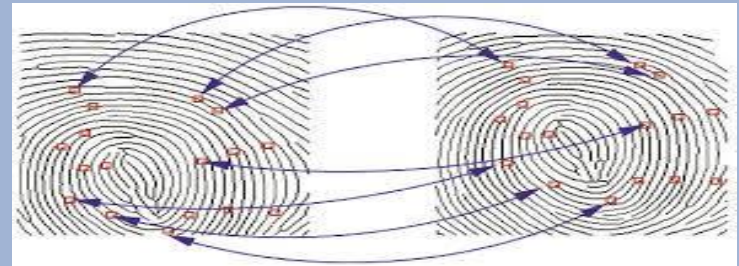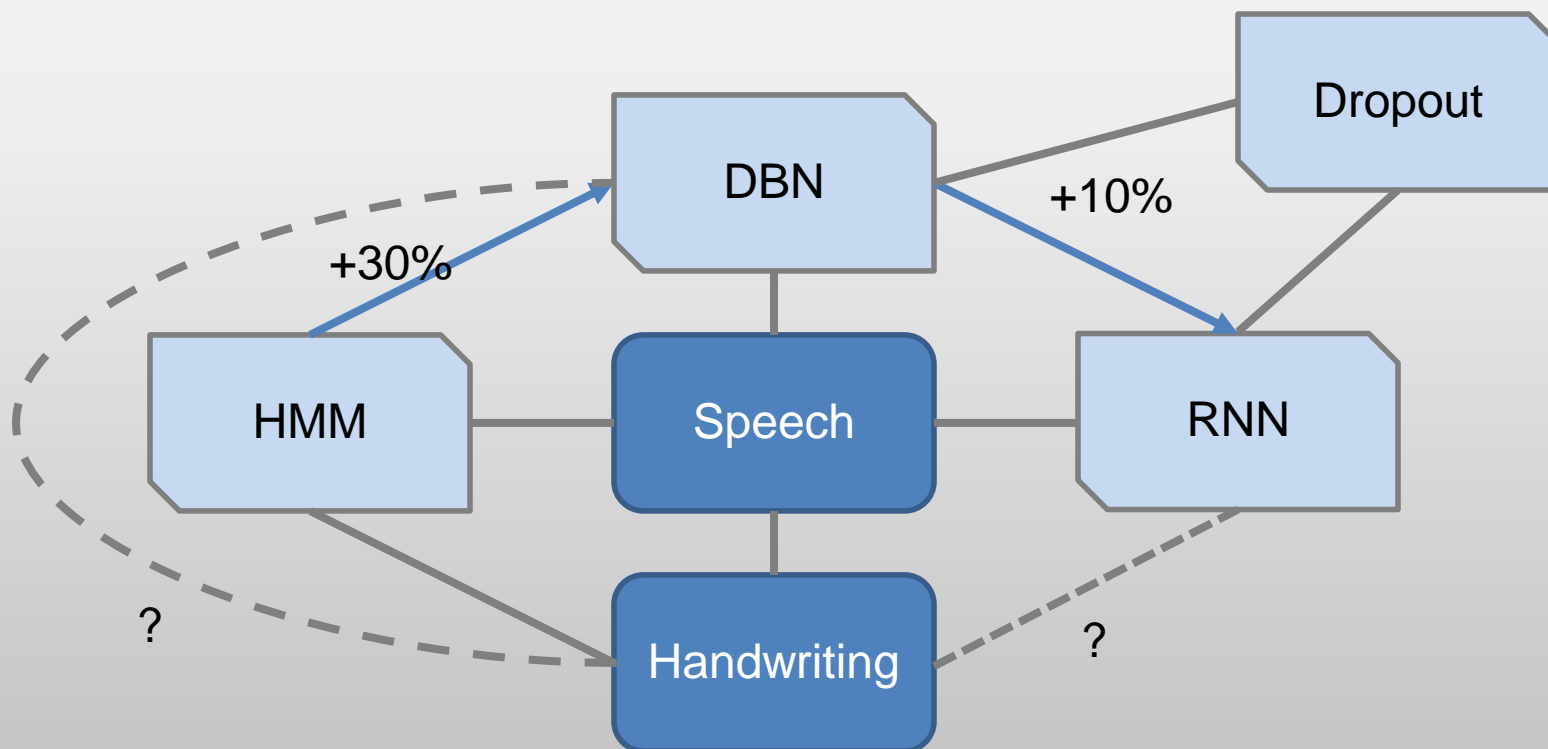
- Legend: -

- Number of lines: 2

# Holistic View



- Online cursive handwriting recognition using speech recognition methods; , John Makhoul, Richard Schwartz, and George Chou  ICASSP 1994
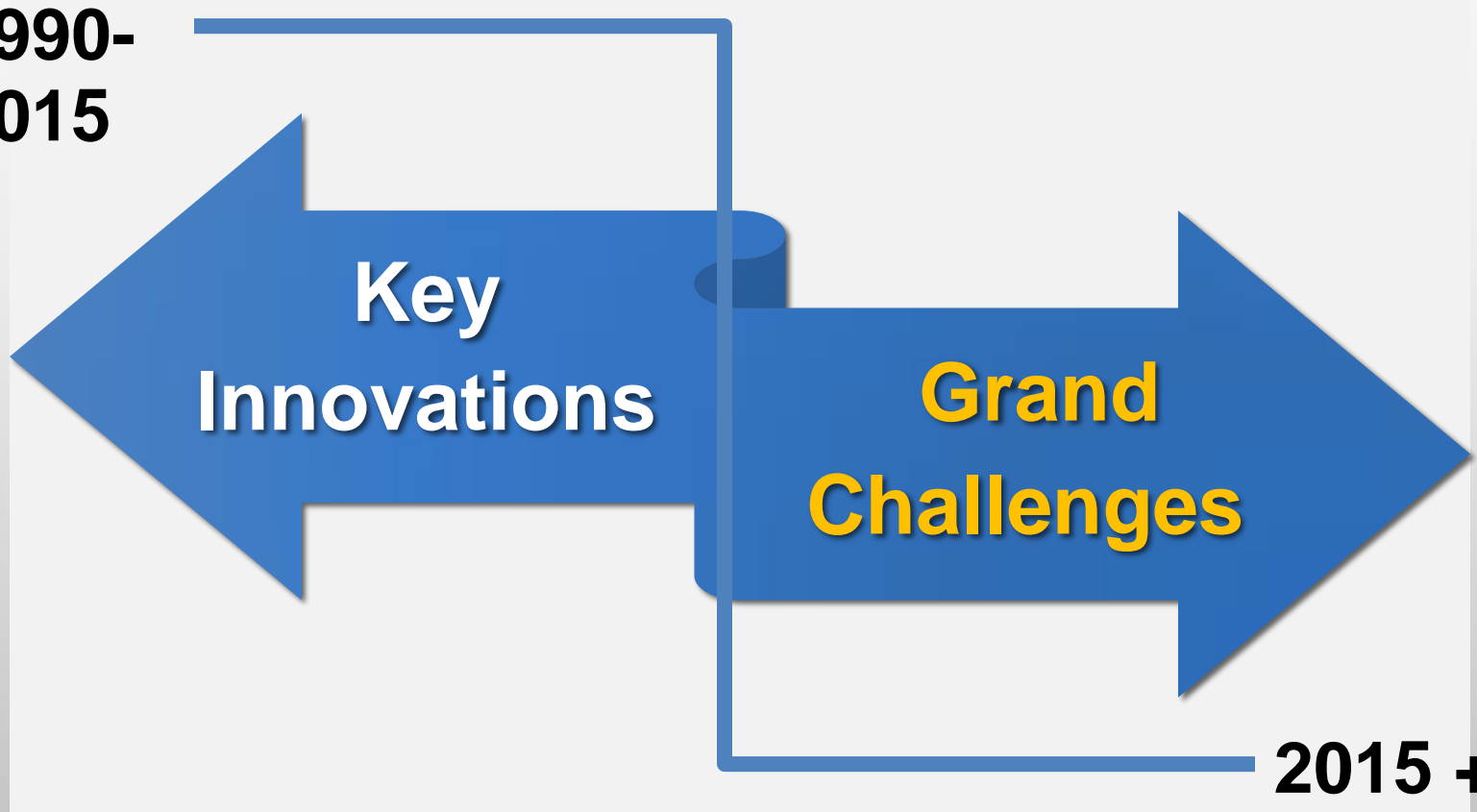
# Accelerated Discovery

**1990-2015**

**Key Innovations**
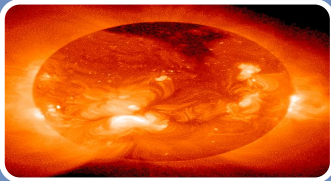
**Grand Challenges**

**2015 +**

# Handwriting Recognition
## Key Innovations


Lexicons


Fusion


Retrieval


Security

# Summary

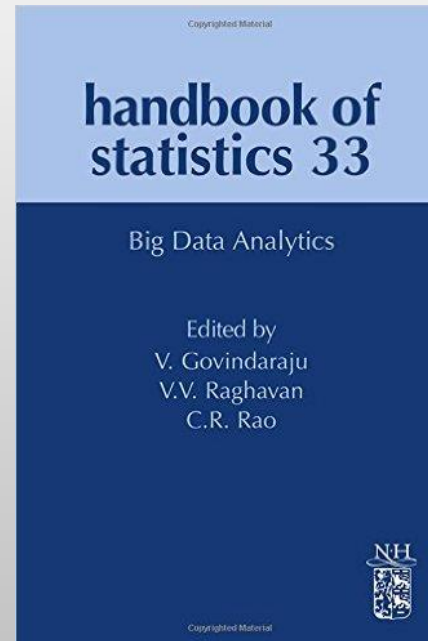| | |
|---|---|
| **Grand Challenges** | • 4Vs of Scientific Big Data<br>• 4 Rs: References, Reinvention, Replicability, Reputation |
| **Grand Opportunities** | • Accelerated Discovery : Supervised linkages, heuristics;<br>• Integrate learning channels |
| **Key Innovations** | • Handwriting Recognition: Lexicons; Fusion; Retrieval; Security |

## Special Thanks to
## All my students and colleagues

## especially to colleagues
## Srirangaraj Setlur and Ifeoma Nwogu



Venu Govinaraju, Ifeoma Nwogu, and Srirangaraj Setlur, "Document Informatics for Scientific Learning and Accelerated Discovery", *Handbook of Statistics* (33): *Big Data Analytics*, pp. 4-28, Elsevier, 2016.

Thank You

venu@cubs.buffalo.edu