# Key Issues in Performance Evaluation for Document Analysis Systems*

## Seventh Mexican Conference on Pattern Recognition
## Mexico City, Mexico
## June 25, 2015

*Professor Daniel Lopresti*
*Computer Science and Engineering*
*Lehigh University*
*Bethlehem, PA, USA*

*\* Copies of these slides will be available on my website next week.*

Key Issues in Performance Evaluation
for Document Analysis Systems

MCPR 2015
June 24-27
Mexico City, Mexico

Lopresti
Slide 1

# Outline

- Motivation / quick overview of document image analysis.
- A simple example of performance evaluation gone awry.
- How do we know when a problem is solved?
- Counting votes – replicating human interpretation.
- A Turing Test-inspired viewpoint.
- Realistic attack models for behavioral biometrics.
- Concluding observations.

*Key Issues in Performance Evaluation*
*for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 2*

# Motivation

Isn't performance evaluation easy?
- We all know accuracy, precision/recall, F-measure, etc.
- Standard datasets and competitions are now common.
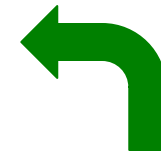
What are the concerns I hear?
- Everyone believes his/her own problem is unique.
- Disconnect between problems and real-world tasks.
- Desperate need to generate publications.
- As a community, we may be too polite.

Why is it so important to do performance evaluation well?
- Measure progress ⇒ prevent wasted effort.
- Scientific respectability.

*Key Issues in Performance Evaluation*
*for Document Analysis Systems*

**MCPR 2015**
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 3*

# Motivation

Typical view of a pattern recognition problem:

*From an advisor or the research literature*

The real world where our solutions must ultimately live:

*Key Issues in Performance Evaluation for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 4*

# Motivation

Pattern recognition techniques are not used in isolation, but rather to solve tasks of interest:

- demands on the degree of automation that is required,
- the minimum acceptable accuracy level,
- and the kinds of errors that can be tolerated.

Important implications for:

- performance evaluation,
- and, ultimately, the success of the system.

*Key Issues in Performance Evaluation*
*for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 5*

# A Quick Overview of Document Image Analysis (DIA) and Optical Character Recognition (OCR)

Many examples I cite draw from this field which presents a rich range of research opportunities. On the other hand, many observations I make will (I hope) generalize.

*Key Issues in Performance Evaluation*
*for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 6*

# Typical DIA Workflow



From "Recognition of textual and graphical patterns" by Josep Lladós, May 2014.

Capture (paper, sketching, fax, electronic formats (PS, PDF, XML, SVG, DXF, ...))

Low-level processing

Analysis & element recognition

Interpretation ↔ domain knowledge

text    layout    graphics

Transfer

Database (DMS)    Re-engineering    Actions    ...

*Key Issues in Performance Evaluation*
*for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 7*

# Document Layout Analysis



Detection of the document logical structure (lines, paragraphs, columns) from the physical structure (image blocks).

Physical Structure

Logical Structure

Document heading

Table heading

Table

Section heading

Text columns

From "Recognition of textual and graphical patterns" by Josep Lladós, May 2014.

*Key Issues in Performance Evaluation for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 8*

# Challenges in Layout Analysis



Big fonts

Backgroud elimination

Logo recognition / advertisements

Non standard text / fonts

*From "Recognition of textual and graphical patterns " by Josep Lladós, May 2014.*

*Key Issues in Performance Evaluation*
*for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 9*

# Components in an OCR Workflow



From "Recognition of textual and graphical patterns " by Josep Lladós, May 2014.

MCPR 2015
June 24-27
Mexico City, Mexico

# Measuring Performance Appropriately: A Simple Worst-Case Example

It's instructive to consider what can go wrong when a standard technique used for performance evaluation is applied without considering the ultimate application.

MCPR 2015
June 24-27
Mexico City, Mexico

# Measuring DIA Performance

construct page grammars sufficiently robust to ignore speckle. This is feasible but tedious. Instead, we filter out all

and tables of contents in technical jour- nals, patent applications, resumes, typed forms with a prespecified layout, sheet

## For the above input, which DIA result is better?

(a) construct page grammors sufficiently robust to ignore spcckle.  This is feasible but tedious.  Instead, we filter out all …

… and tobles of contents in technical journals, patent appllcations, resumes, typed fonns with a prespecified 1ayout, sheet …

*OCR errors*

(b) construct page grammars sufficiently and tables of contents in technical jour- robust to ignore speckle.  This is feasible nals, patent applications, resumes, typed but tedious.  Instead, we filter out all forms with a prespecified layout, sheet …

*Missed columns*

*Key Issues in Performance Evaluation*
*for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 12*

# Measuring DIA Performance

Which DIA result is better?

It depends on the application!

| DIA | → | OCR errors | → | Text to Speech |
| DIA | → | Column errors (crossed out) | → | Text to Speech |

| DIA | → | OCR errors (crossed out) | → | Bag of Words IR |
| DIA | → | Column errors | → | Bag of Words IR |

*Key Issues in Performance Evaluation*
*for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 13*

# String Edit Distance

Edit distance is standard measure used for OCR accuracy.

**The Edit Distance Problem.**

Given two sequences, find the optimal series of deletions, insertions, and substitutions to transform one into the other.

**Input:** Two sequences:

$$v = v_1\ v_2 \ldots v_m \quad \text{and} \quad w = w_1\ w_2 \ldots w_n$$

**Output:** An optimal series of basic editing operations:

$$e_1,\ e_2,\ \ldots,\ e_t$$

such then, when applied to one of the sequences, say $v$, it is transformed into the other sequence, $w$.

Here "optimal" can mean any of a number of things, including "fewest" or "lowest- / highest-cost."

*Key Issues in Performance Evaluation
for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti
Slide 14*

# String Edit Distance

Given two sequences $v$ and $w$, consider what is required to compute optimal distance between prefixes $v[1..i]$ and $w[1..j]$.  There are three possible cases:

| deletion | insertion | substitution or match |
|----------|-----------|------------------------|

| optimal distance for $v[1..i\text{-}1]$ and $w[1..j]$ | $v[i]$ <br> ↕ <br> - | optimal distance for $v[1..i]$ and $w[1..j\text{-}1]$ | - <br> ↕ <br> $w[j]$ | optimal distance for $v[1..i\text{-}1]$ and $w[1..j\text{-}1]$ | $v[i]$ <br> ↕ *or* ↕ <br> $w[j]$ |

<p align="center"><b>I      II      III</b></p>

If we already have solutions for all shorter prefixes, we can compute the distance for $v[1..i]$ and $w[1..j]$.

*Key Issues in Performance Evaluation for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 15*

# String Edit Distance

Conceptually, this might look something like this:

$$\text{optimal distance at } v[1..i] \text{ and } w[1..j] = \min \begin{cases} \text{optimal distance at } v[1..i-1] \text{ and } w[1..j] \\ + \\ \text{cost of deleting } v[i] \\ \\ \text{optimal distance at } v[1..i] \text{ and } w[1..j-1] \\ + \\ \text{cost of inserting } w[j] \\ \\ \text{optimal distance at } v[1..i-1] \text{ and } w[1..j-1] \\ + \\ \text{cost of matching } v[i] \text{ and } w[j] \end{cases}$$

We assume deletions, insertions, and mismatches have positive cost, while matches have zero or negative cost.

*Key Issues in Performance Evaluation for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 16*

# Measuring DIA Performance

construct page grammars sufficiently robust to ignore speckle.  This is feasible but tedious.  Instead, we filter out all

...

and tables of contents in technical jour-nals, patent applications, resumes, typed forms with a prespecified layout, sheet

construct page grammors sufficiently robust to ignore spcckle.  This is feasible but tedious.  Instead, we filter out all

...

and tobles of contents in technical jour-nals, patent appllcations, resumes, typed fonns with a prespecified 1ayout, sheet

*Small edit distance*

*Ground truth*

*Large edit distance*

construct page grammars sufficiently and tables of contents in technical jour-robust to ignore speckle.  This is feasible nals, patent applications, resumes, typed but tedious.  Instead, we filter out all forms with a prespecified layout, sheet

It is vital to match your performance measure to your target application.

*Key Issues in Performance Evaluation
for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti
Slide 17*

Just now "in the weeds."     Turn to 30,000 foot view.



Specific details of performance evaluation
➔ When is a problem solved?

Performance evaluation confirms when we have advanced state of the art and, ultimately, when a problem is solved.

*Key Issues in Performance Evaluation*
*for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 18*

# What's the challenge?

We define our open problems as automating a task:  this is quite different from math, physics, theoretical CS, etc.

Some ways of measuring success:

- Accuracy of new algorithm (vs. previous methods).
- Current degree of community interest (publishability).
- Economic considerations (net payoff for using method).
- Distinguishability of algorithm from human result.

"When is a Problem Solved?," D. Lopresti and G. Nagy, Proceedings of the Eleventh International Conference on Document Analysis and Recognition (ICDAR 2011), September 2011, Beijing, China, pp. 32-36.

*Key Issues in Performance Evaluation for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 19*

# What's the challenge?

When is a problem solved?

This seems like a simple, basic question.

It also seems like an important question.

But it's not clear we know how to answer it ...

"When is a Problem Solved?," D. Lopresti and G. Nagy, Proceedings of the Eleventh International Conference on Document Analysis and Recognition (ICDAR 2011), September 2011, Beijing, China, pp. 32-36.

*Key Issues in Performance Evaluation for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 20*

# Viewpoint #1

The endless pursuit of perfection:

"A problem is solved if there is a method which has been widely publicized and documented and freely available to the community which achieves 100% accuracy on within-spec inputs it receives."

"When is a Problem Solved?," D. Lopresti and G. Nagy, Proceedings of the Eleventh International Conference on Document Analysis and Recognition (ICDAR 2011), September 2011, Beijing, China, pp. 32-36.

*Key Issues in Performance Evaluation*
*for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 21*

# Viewpoint #2

As good as it gets:

"A problem is solved if there is a method which has been widely publicized and documented and freely available to the community which performs better than any other method, and which cannot be further improved without investing excessive resources."

"When is a Problem Solved?," D. Lopresti and G. Nagy, Proceedings of the Eleventh International Conference on Document Analysis and Recognition (ICDAR 2011), September 2011, Beijing, China, pp. 32-36.

Key Issues in Performance Evaluation
for Document Analysis Systems

MCPR 2015
June 24-27
Mexico City, Mexico

Lopresti
Slide 22

# Viewpoint #3

Good enough to get the job done:

> "A problem is solved if there is a method which has been widely publicized and documented and freely available to the community which cannot be replaced with any other method to improve the end-to-end performance of a specific application of interest."

"When is a Problem Solved?," D. Lopresti and G. Nagy, Proceedings of the Eleventh International Conference on Document Analysis and Recognition (ICDAR 2011), September 2011, Beijing, China, pp. 32-36.

*Key Issues in Performance Evaluation for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti
Slide 23*

# Viewpoint #4

Pure pragmatism:

> "A problem is solved when it is no longer possible to get a paper published on the topic (or, alternatively, to raise research funding to study the question)."

"When is a Problem Solved?," D. Lopresti and G. Nagy, Proceedings of the Eleventh International Conference on Document Analysis and Recognition (ICDAR 2011), September 2011, Beijing, China, pp. 32-36.

*Key Issues in Performance Evaluation for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti
Slide 24*

# Viewpoint #5

The Turing Test:

> "A problem is solved if there is a method which has been widely publicized and documented and freely available to the community which generates output for a given input that a human judge cannot reliably distinguish from the output of a human expert."

"When is a Problem Solved?," D. Lopresti and G. Nagy, Proceedings of the Eleventh International Conference on Document Analysis and Recognition (ICDAR 2011), September 2011, Beijing, China, pp. 32-36.

Key Issues in Performance Evaluation
for Document Analysis Systems

MCPR 2015
June 24-27
Mexico City, Mexico

Lopresti
Slide 25

# Show of Hands

Which viewpoint(s) do you agree with?

- The endless pursuit of perfection.

- As good as it gets.

- Good enough to get the job done.

- Pure pragmatism.

- The Turing Test.

"When is a Problem Solved?," D. Lopresti and G. Nagy, Proceedings of the Eleventh International Conference on Document Analysis and Recognition (ICDAR 2011), September 2011, Beijing, China, pp. 32-36.

# Key Considerations

Let's also keep in mind the following important points:

- Populations vs. samples:  performance figures like error, reject, or retrieval rates are of interest only with regard to populations rather than specific samples.

- Algorithms, heuristics, and implementations:  most of pattern recognition is built on heuristics rather than algorithms, although the latter term is applied to both. To be a solution, an algorithm must be implementable.

- Desirable criteria:  solutions should be invariant to 90° rotation, modest differences in resolution, remapping RGB/gray values, jitters in threshold settings, etc.

"When is a Problem Solved?," D. Lopresti and G. Nagy, Proceedings of the Eleventh International Conference on Document Analysis and Recognition (ICDAR 2011), September 2011, Beijing, China, pp. 32-36.

Key Issues in Performance Evaluation
for Document Analysis Systems

MCPR 2015
June 24-27
Mexico City, Mexico

Lopresti
Slide 27

# A Simple Yet Vexing Case Study: Counting Votes Recorded on Paper

Topic of current interest where the legal need to respect voter intent transforms a seemingly trivial pattern recognition problem into much more complex task.

*Key Issues in Performance Evaluation*
*for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 28*

# Counting Votes Not So Easy



**INSTRUCTIONS TO VOTERS**
To vote, completely fill in the oval(s) next to your choice(s) like this:

Is this a legal vote?

- Courts would probably say so ...
- ... but op-scan readers might not count it.

*Increasing demands that machine's interpretation match a human's.*

*Key Issues in Performance Evaluation*
*for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 29*

# Research Questions

Issues that arise from using paper ballots in elections:

- Accurate interpretation of marginal markings.
- Human cost, error rate, and bias in performing manual recounts.
- Failure modes in ballot imaging (e.g., paper jams).
- Systematic errors due to ballot layout (one candidate may be disadvantaged over another based on physical location on page).

Also keep in mind:

- U.S. elections can be complex (10's to 100's of choices).
- Impact of "voter error" (e.g., improper markings, erasures).
- Potential for traditional ballot-box stuffing.
- Computer hackers attempting to manipulate the vote.

*Key Issues in Performance Evaluation for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 30*

# Why isn't this a solved problem?

Students have been taking standardized tests using op-scan answer sheets for decades …





- While accuracy rates are very high, problems do occur.
- Compared to voters, students are a much more homogeneous (and well-educated) population.
- Standardized testing is NOT anonymous.  Students can (and do) complain when they receive a lower score than they expect.

*Key Issues in Performance Evaluation*
*for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 31*

# Connection to Forms Processing

Similarities to forms processing, but also some key differences:

- Much broader range of users (education level, literacy, etc.) than for traditional forms applications.

- Ballots must preserve a voter's anonymity.

- Demand to count votes and report results quickly.

- Elections are held infrequently, so voting equipment sits unused for long periods in storage.

- Poll workers often lack technical expertise.

- Maintaining chain-of-custody is a critical security requirement.

- No *financial* interest in making sure votes are counted accurately, but there is tremendous *public* interest.

*Key Issues in Performance Evaluation for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 32*

# Counting Votes Not So Easy

Real ballot from an election in California:



One of these votes was counted correctly by the op-scan equipment, the other was not.

Note: this does <u>not</u> mean voting on paper ballots is bad, just (1) manual audits should be mandatory, and (2) more research is needed.

"Improving California's 1% Manual Tally Procedure," Joseph Lorenzo Hall, UC Berkeley School of Information, EVT Workshop 2008.

*Key Issues in Performance Evaluation for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 33*

# Whole-Ballot Recognition



**Stray mark?**

**Valid vote?**

BASEBALL HALL OF FAME
(vote for no more than 5)

Ty Cobb
Rogers Hornsby
Walter Johnson
Nap Lajoie
Christy Mathewson
Babe Ruth
Tris Speaker
Honus Wagner
Cy Young

Can we capture voter intent via style-based techniques?

*Key Issues in Performance Evaluation*
*for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 34*

# Style-Based Mark Recognition

**Traditional Forms Processing**

**Style-Base Ballot Mark Recognition**

*Can the system interpret the voter's intent?* (If a human judge would interpret a marking as an intended vote, then the voting machine should do the same.)

Can fail to record some votes simply because they do not satisfy an arbitrary criterion (e.g., a fixed threshold on the number of black pixels).

Assume a voter is self-consistent when marking his/her ballot.

Create a style-based classifier from a set of style-specialized classifiers to improve recognition accuracy.

**Limiting**

**Promising**

"Style-Based Ballot Mark Recognition," P. Xiu, D. Lopresti, H. Baird, G. Nagy, and E. Barney Smith, *Proceedings of the Tenth International Conference on Document Analysis and Recognition*, July 2009, Barcelona, Spain, pp. 216-220.

*Key Issues in Performance Evaluation for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 35*

# Challenging Cases

**A Style-Consistent Field**

**Human Interpretation**

*VN*NN, NVNN, NNVN, NVNN, NNNV

**(V for Vote, N for Non-vote)**

**Variations in Marking Style**

**Check, ex, and filled marks (left)
vs. noisy non-votes (right)**

"Style-Based Ballot Mark Recognition," P. Xiu, D. Lopresti, H. Baird, G. Nagy, and E. Barney Smith, *Proceedings of the Tenth International Conference on Document Analysis and Recognition*, July 2009, Barcelona, Spain, pp. 216-220.

# System Design



"Style-Based Ballot Mark Recognition," P. Xiu, D. Lopresti, H. Baird, G. Nagy, and E. Barney Smith, *Proceedings of the Tenth International Conference on Document Analysis and Recognition*, July 2009, Barcelona, Spain, pp. 216-220.

# Style-Based Performance

## Table 3. Target-level error rates (top) and field-level error rates (bottom).

| Sample Set | Classifier | | | | | |
| --- | --- | --- | --- | --- | --- | --- |
| | Check | Ex | Filled | Blend | Separate | Style-based |
| Check | 2.36% | 7.46% | 25.00% | 1.97% | 4.35% | 2.78% |
| Ex | 0.40% | 0.34% | 16.16% | 0.40% | 0.40% | 0.35% |
| Filled | 2.75% | 2.38% | 1.10% | 2.75% | 2.50% | 1.09% |
| Average | 1.84% | 3.39% | 14.09% | 1.70% | 2.42% | 1.41% |

| Sample Set | Classifier | | | | | |
| --- | --- | --- | --- | --- | --- | --- |
| | Check | Ex | Filled | Blend | Separate | Style-based |
| Check | 38.30% | 83.25% | 100.00% | 33.43% | 61.08% | 42.85% |
| Ex | 7.77% | 6.70% | 99.30% | 7.77% | 7.77% | 6.75% |
| Filled | 53.18% | 46.07% | 20.75% | 53.18% | 48.55% | 20.63% |
| Average | 33.08% | 45.34% | 73.35% | 31.46% | 39.13% | 23.41% |

"Style-Based Ballot Mark Recognition," P. Xiu, D. Lopresti, H. Baird, G. Nagy, and E. Barney Smith, *Proceedings of the Tenth International Conference on Document Analysis and Recognition*, July 2009, Barcelona, Spain, pp. 216-220.

*Key Issues in Performance Evaluation for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 38*

# BallotGen Mark Synthesis

*Key Issues in Performance Evaluation*
*for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 39*

# A Bit of Good Luck

But what we'd like to have is ballots from a real election. Even better, the ballots would be from an important election where the voter markings present serious pattern recognition challenges.

Extremely close U.S. Senate race in State of Minnesota: six days after election, unofficial results showed Republican Norm Coleman leading Democratic challenger Al Franken by 206 votes out of nearly 3 million cast, a difference of less than 0.01%.

"Document Analysis Issues in Reading Optical Scan Ballots," D. Lopresti, G. Nagy, and E. Barney Smith, *Proceedings of the Ninth IAPR International Workshop on Document Analysis Systems*, June 2010, Boston, MA, pp. 105-112.

*Key Issues in Performance Evaluation for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 40*

# A Bit of Good Luck

- Minnesota uses op-scan ballots.

- Closeness of election triggers a manual recount.

- Both sides are allowed to challenge validity of "questionable" ballots.

- Openness laws make challenged ballots a matter of public record.

- Ballot images made available on MN public radio website.

- PDF files contain 300 dpi TIF images!

http://minnesota.publicradio.org/features/2008/11/19_challenged_ballots/

*Key Issues in Performance Evaluation for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 41*

# Minnesota Statutes

Remember that the guiding principle is *voter intent*. Here are a few key points to keep in mind when interpreting ballot markings:

- "A ballot shall not be rejected for a technical error that does not make it impossible to determine the voter's intent."

- "If a mark (X) is made out of its proper place, but so near a name or space as to indicate clearly the voter's intent, the vote shall be counted."

- "Misspelling or abbreviations of the names of write-in candidates shall be disregarded if the individual for whom the vote was intended can be clearly ascertained from the ballot."

https://www.revisor.mn.gov/statutes/?id=204C.22

MCPR 2015
June 24-27
Mexico City, Mexico

# Minnesota Statutes

... and ...

- "If a voter uniformly uses a mark other than (X) which clearly indicates an intent to mark a name or to mark yes or no on a question, and the voter does not use (X) anywhere else on the ballot, a vote shall be counted for each candidate or response to a question marked.

- If a voter uses two or more distinct marks, such as (X) and some other mark, a vote shall be counted for each candidate or response to a question marked, unless the ballot is marked by distinguishing characteristics that make the entire ballot defective ..."

https://www.revisor.mn.gov/statutes/?id=204C.22

MCPR 2015
June 24-27
Mexico City, Mexico

# Minnesota Statutes

… and …

- "If the names of two candidates have been marked, and an attempt has been made to erase or obliterate one of the marks, a vote shall be counted for the remaining marked candidate."

- "A ballot shall not be rejected merely because it is slightly soiled or defaced."

- "If a ballot is marked by distinguishing characteristics in a manner making it evident that the voter intended to identify the ballot, the entire ballot is defective."

https://www.revisor.mn.gov/statutes/?id=204C.22

Goal here is to prevent coercion or vote selling.

*Key Issues in Performance Evaluation*
*for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 44*

# Challenge: you be the judge



Who gets vote?  Public opinion:

- Norm Coleman:  63% (7,626 votes)
- Al Franken:  4% (474 votes)
- Nobody:  33% (4,050 votes)

*Key Issues in Performance Evaluation*
*for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 45*

# Challenge: you be the judge



Vote for Franken?  Public opinion:
- Yes:  92% (11,069 votes)
- No:  8% (1,012 votes)

*Key Issues in Performance Evaluation*
*for Document Analysis Systems*

**MCPR 2015**
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 46*

# Challenge:  you be the judge



Vote for Franken?  Public opinion:

- Yes:  96% (11,250 votes)
- No:  4% (452 votes)

MCPR 2015
June 24-27
Mexico City, Mexico

# Challenge: you be the judge



Vote for Coleman?  Public opinion:

- Yes:  54% (6,080 votes)
- No:  46% (5,203 votes)

*Key Issues in Performance Evaluation*
*for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 48*

# MN Challenged Ballot Collection

How the ballot collection was generated and harvested:

- Ballots photocopied and originals stored in a secure location.

- Copies scanned to PDF using auto-feeder flatbed scanner.

- Ballot was two-sided, with both sides scanned simultaneously.

- I wrote a simple web "crawler" that automatically downloaded all the files and extracted TIF images from PDF.

- A total of 6,737 ballots in the set.

- Examination of the TIF suggests that ballots were scanned at 300 dpi bitonal, and that lossy compression was never used.

- Hence, they form an ideal dataset for research purposes.

*Key Issues in Performance Evaluation*
*for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 49*

# Minnesota Ballot Front and Back

MCPR 2015
June 24-27
Mexico City, Mexico

# Sloppy-But-Valid Marks



*Key Issues in Performance Evaluation*
*for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 51*

# Non-Conforming Marking Styles

*Key Issues in Performance Evaluation*
*for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 52*

# Attempts to Cancel a Vote

*Key Issues in Performance Evaluation*
*for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 53*

# Votes that Look Cancelled

# Stray Marks and Bleedthrough

*Key Issues in Performance Evaluation*
*for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 55*

# Invalidating Markings

*Key Issues in Performance Evaluation*
*for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 56*

# An Example from Mexico

And such issues are not limited to U.S. elections …

*Key Issues in Performance Evaluation for Document Analysis Systems*

*Lopresti*
*Slide 57*

MCPR 2015
June 24-27
Mexico City, Mexico

# Why isn't this an easy problem?

After all, ballots are just a simple type of form.  We must read votes correctly, but we aren't expected to recognize write-ins.

Can't we just push up reject rate until accuracy reaches 100%?

> Remember, we can't change rules in ways that violate the law. VOTER INTENT is the definition we must always follow.

To do this right, we must be prepared to:

- Reject any ballot that may contain "identifying marks."

- Recognize intent when mark is atypical or far from target.

- Accurately identify when a vote has been cancelled.

*Key Issues in Performance Evaluation*
*for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 58*

# Status

- Ground truth collected from 8 test subjects, 980 ballot sides.

- All 6,737 ballots now online on DAE server (see URL below for more details on the server and its capabilities).



http://dae.cse.lehigh.edu/DAE/

*Key Issues in Performance Evaluation for Document Analysis Systems*

**MCPR 2015**
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 59*

# Adapting the Turing Test for Declaring a Problem Solved

An interesting thought experiment, given the demand for algorithms that can perform at human levels when users are free to act in ways that confound the system.

*Key Issues in Performance Evaluation*
*for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 60*

# Recall from Earlier ...

The Turing Test:

> "A problem is solved if there is a method which has been widely publicized and documented and freely available to the community which generates output for a given input that a human judge cannot reliably distinguish from the output of a human expert."

Differs significantly from employing ground-truth provided by a human expert in advance.

"Adapting the Turing Test for Declaring Document Analysis Problems Solved," D. Lopresti and G. Nagy, Proceedings of the Tenth IAPR International Workshop on Document Analysis Systems (DAS 2012), March 2012, Gold Coast, Australia, 5 pages.

Key Issues in Performance Evaluation
for Document Analysis Systems

MCPR 2015
June 24-27
Mexico City, Mexico

Lopresti
Slide 61

# The Imitation Game

1. *The Imitation Game.*

I PROPOSE to consider the question, 'Can machines think?'
This should begin with definitions of the meaning of the terms
'machine' and 'think'. The definitions might be framed so as to
reflect so far as possible the normal use of the words, but this
attitude is dangerous. If the meaning of the words 'machine'
and 'think' are to be found by examining how they are commonly
used it is difficult to escape the conclusion that the meaning
and the answer to the question, 'Can machines think?' is to be
sought in a statistical survey such as a Gallup poll. But this is
absurd. Instead of attempting such a definition I shall replace the
question by another, which is closely related to it and is expressed
in relatively unambiguous words.

The new form of the problem can be described in terms of
a game which we call the 'imitation game'. It is played with
three people, a man (A), a woman (B), and an interrogator (C) who
may be of either sex. The interrogator stays in a room apart
from the other two. The object of the game for the interrogator
is to determine which of the other two is the man and which is
the woman. He knows them by labels X and Y, and at the end
of the game he says either 'X is A and Y is B' or 'X is B and Y
is A'. The interrogator is allowed to put questions to A and B
thus:

C: Will X please tell me the length of his or her hair?
Now suppose X is actually A, then A must answer. It is A's

28        433

---

1. *The Imitation Game.*

I PROPOSE to consider the question, 'Can machines think?'
This should begin with definitions of the meaning of the terms
'machine' and 'think'. The definitions might be framed so as to

---

The new form of the problem can be described in terms of
a game which we call the 'imitation game'. It is played with
three people, a man (A), a woman (B), and an interrogator (C) who
may be of either sex. The interrogator stays in a room apart
from the other two. The object of the game for the interrogator
is to determine which of the other two is the man and which is
the woman. He knows them by labels X and Y, and at the end

---

We now ask the question, 'What will happen when a machine
takes the part of A in this game?' Will the interrogator decide
wrongly as often when the game is played like this as he does
when the game is played between a man and a woman? These
questions replace our original, 'Can machines think?'

A. M. Turing, "Computing Machinery and Intelligence,"
*Mind*, vol. 59, no. 236, October 1950, pp. 433-460.

"Adapting the Turing Test for Declaring Document Analysis Problems Solved," D. Lopresti and G. Nagy, Proceedings of the Tenth IAPR
International Workshop on Document Analysis Systems (DAS 2012), March 2012, Gold Coast, Australia, 5 pages.

Key Issues in Performance Evaluation
for Document Analysis Systems

MCPR 2015
June 24-27
Mexico City, Mexico

Lopresti
Slide 62

# The Turing Test



(A) Man pretending to be woman

(B) Woman trying to help interrogator

(C) Interrogator trying to make right guess

SuccessRate$_1$

(A) Machine pretending to be woman

(B) Woman trying to help interrogator

(C) Interrogator trying to make right guess

SuccessRate$_2$

Is SuccessRate$_2$ ≈ SuccessRate$_1$ ?

"Adapting the Turing Test for Declaring Document Analysis Problems Solved," D. Lopresti and G. Nagy, Proceedings of the Tenth IAPR International Workshop on Document Analysis Systems (DAS 2012), March 2012, Gold Coast, Australia, 5 pages.

*Key Issues in Performance Evaluation for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 63*

# The Turing Test

The Turing Test is an elegantly simple idea, so it should be simple to implement, right?



- (A) Machine performing some task
- (B) Human performing same task
- (C) Interrogator trying to make right guess

Is SuccessRate no better than random chance ?

- Note this differs from Turing's original formulation.

- When considering a real implementation, other, more serious complications arise.

"Adapting the Turing Test for Declaring Document Analysis Problems Solved," D. Lopresti and G. Nagy, Proceedings of the Tenth IAPR International Workshop on Document Analysis Systems (DAS 2012), March 2012, Gold Coast, Australia, 5 pages.

*Key Issues in Performance Evaluation*
*for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 64*

# Long Bet*



"By 2029 no computer – or 'machine intelligence' – will have passed the Turing Test."

PREDICTOR:
    Mitchell Kapor
CHALLENGER:
    Ray Kurzweil
STAKES:  $20,000

"Adapting the Turing Test for Declaring Document Analysis Problems Solved," D. Lopresti and G. Nagy, Proceedings of the Tenth IAPR International Workshop on Document Analysis Systems (DAS 2012), March 2012, Gold Coast, Australia, 5 pages.

*Key Issues in Performance Evaluation for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 65*

# Long Bet Rules

Turing was nonspecific about how to administer his Test, but concreteness is needed when $20,000 is at stake.

- Each of three Turing Test judges is to conduct an online interview ("chat") with each of four human players as well as the machine for two hours.

- At the end of these interviews, the judges indicate whether or not each candidate is human and also rank them from "least human" to "most human."

- The machine is said to pass the Turing Test if it fools two or more judges and if its median rank is equal to or greater than at least two of the human players.

"Adapting the Turing Test for Declaring Document Analysis Problems Solved," D. Lopresti and G. Nagy, Proceedings of the Tenth IAPR International Workshop on Document Analysis Systems (DAS 2012), March 2012, Gold Coast, Australia, 5 pages.
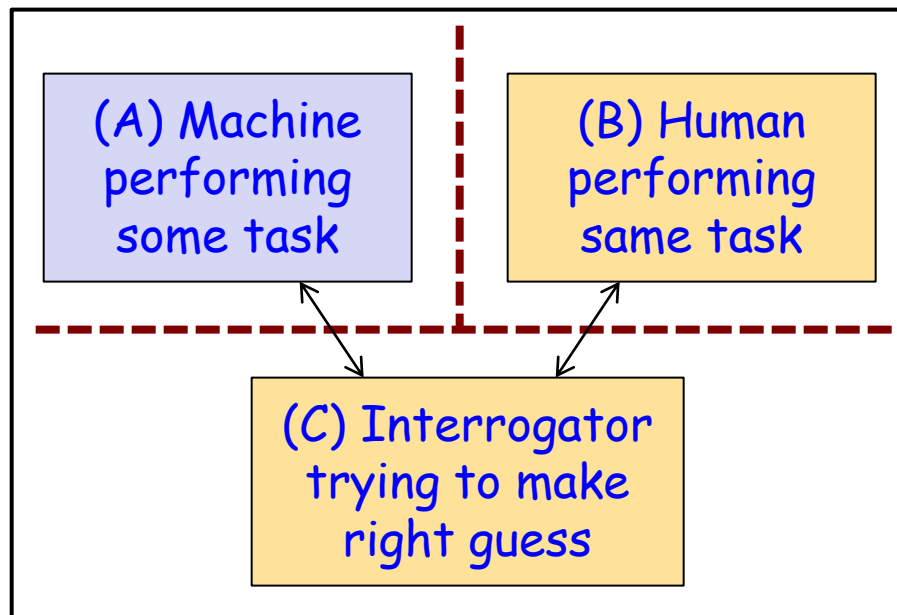
*Key Issues in Performance Evaluation for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 66*

# Adapting the Turing Test

The Long Bet is a one-time event with a significant amount of prize money involved.  As a result, it makes sense to employ a heavy-weight protocol for the test.

How can the Turing Test be applied in document analysis?

- What are the essential qualities to preserve?
- What can be dispensed with, or at least simplified?
- When implemented, how would the test "look"?
- When might such a test be appropriate?

"Adapting the Turing Test for Declaring Document Analysis Problems Solved," D. Lopresti and G. Nagy, Proceedings of the Tenth IAPR International Workshop on Document Analysis Systems (DAS 2012), March 2012, Gold Coast, Australia, 5 pages.

*Key Issues in Performance Evaluation*
*for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 67*

# Properties to Preserve #1

Human judgment is applied to determine a simple machine/human distinction and nothing more complex than this. Automated evaluation (i.e., a computation to determine how "similar" a machine output is to some predefined human "ground truth") is ruled out.



"Adapting the Turing Test for Declaring Document Analysis Problems Solved," D. Lopresti and G. Nagy, Proceedings of the Tenth IAPR International Workshop on Document Analysis Systems (DAS 2012), March 2012, Gold Coast, Australia, 5 pages.

MCPR 2015
June 24-27
Mexico City, Mexico

# Properties to Preserve #2

A judge may ask any number of questions before making a determination. A "question" here is a challenge that requires a response from the player. For document analysis applications, this will normally consist of a page image to be processed in some way.



"Adapting the Turing Test for Declaring Document Analysis Problems Solved," D. Lopresti and G. Nagy, Proceedings of the Tenth IAPR International Workshop on Document Analysis Systems (DAS 2012), March 2012, Gold Coast, Australia, 5 pages.

Key Issues in Performance Evaluation
for Document Analysis Systems

MCPR 2015
June 24-27
Mexico City, Mexico

Lopresti
Slide 69

# Properties to Preserve #3

The judge decides which questions to use, and is free to conduct the questioning of the players without constraint on the choice, sequence, and number of questions.



"Adapting the Turing Test for Declaring Document Analysis Problems Solved," D. Lopresti and G. Nagy, Proceedings of the Tenth IAPR International Workshop on Document Analysis Systems (DAS 2012), March 2012, Gold Coast, Australia, 5 pages.
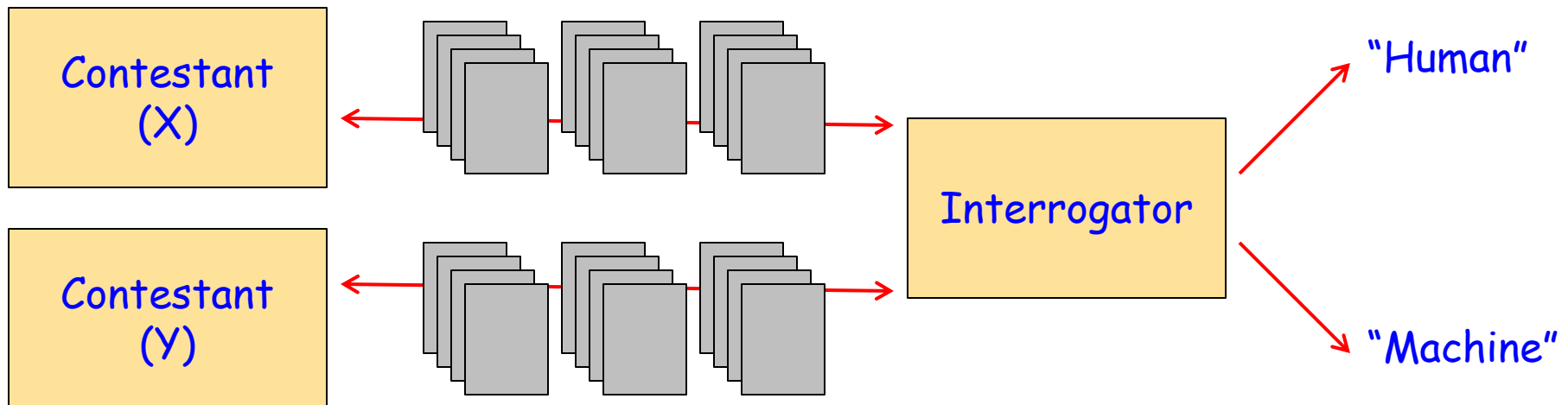
*Key Issues in Performance Evaluation*
*for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 70*

# Properties to Preserve #4

A series of such evaluations, with anyone being allowed to volunteer to serve as judge or as the human player, is conducted before declaring a problem "solved" (if/when the success rates of the best-performing judges are statistically no better than random).



"Adapting the Turing Test for Declaring Document Analysis Problems Solved," D. Lopresti and G. Nagy, Proceedings of the Tenth IAPR International Workshop on Document Analysis Systems (DAS 2012), March 2012, Gold Coast, Australia, 5 pages.
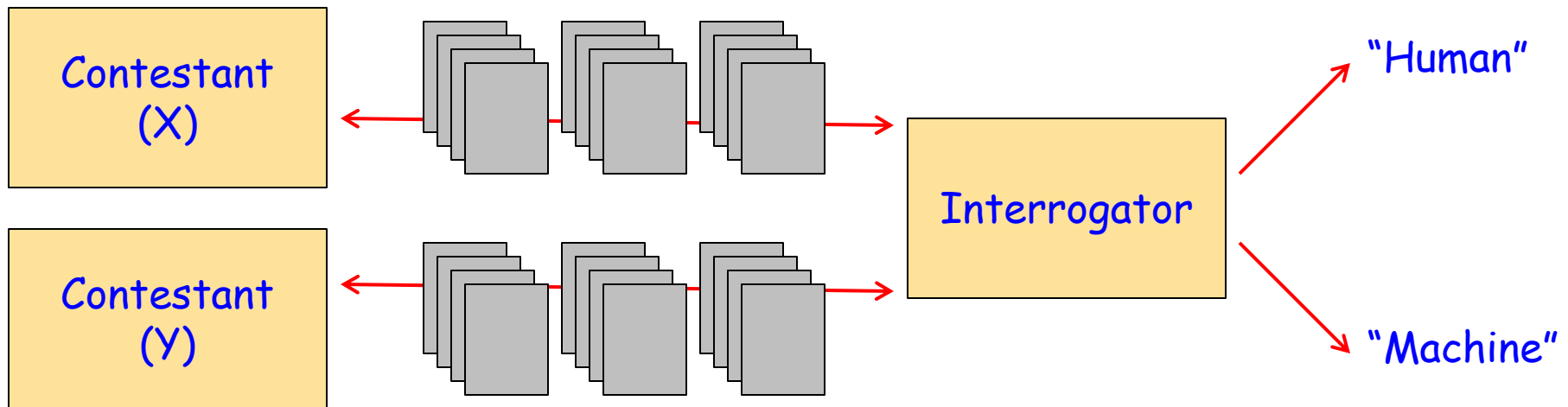
Key Issues in Performance Evaluation
for Document Analysis Systems

MCPR 2015
June 24-27
Mexico City, Mexico

Lopresti
Slide 71

# Properties to Adapt

Some aspects of Turing's original Test must be updated:

- The judge and players do not interact via a natural language question-and-answer process. Instead, they employ a graphical user interface which supports the upload of image files and visual inspection of results.

- The domain of discourse is no longer open-ended. Note that this replaces Turing's original question "Can machines think?" with our "Is this problem solved?"

"Adapting the Turing Test for Declaring Document Analysis Problems Solved," D. Lopresti and G. Nagy, Proceedings of the Tenth IAPR International Workshop on Document Analysis Systems (DAS 2012), March 2012, Gold Coast, Australia, 5 pages.
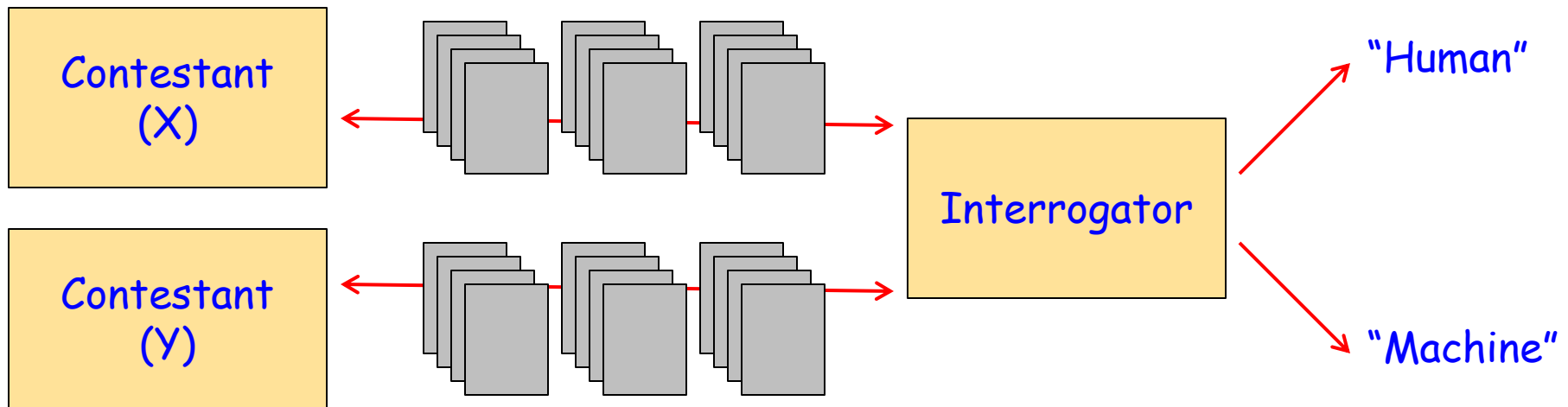
# GUI from Judge's Perspective

Task is:  Logo Detection        Current Challenge is #12

**Pre-defined Challenge Library**



**Create New Challenge**

File name        Upload

Submit to Player A        Submit to Player B

Responses



Determination:        A human, B machine        A machine, B human

"Adapting the Turing Test for Declaring Document Analysis Problems Solved," D. Lopresti and G. Nagy, Proceedings of the Tenth IAPR International Workshop on Document Analysis Systems (DAS 2012), March 2012, Gold Coast, Australia, 5 pages.
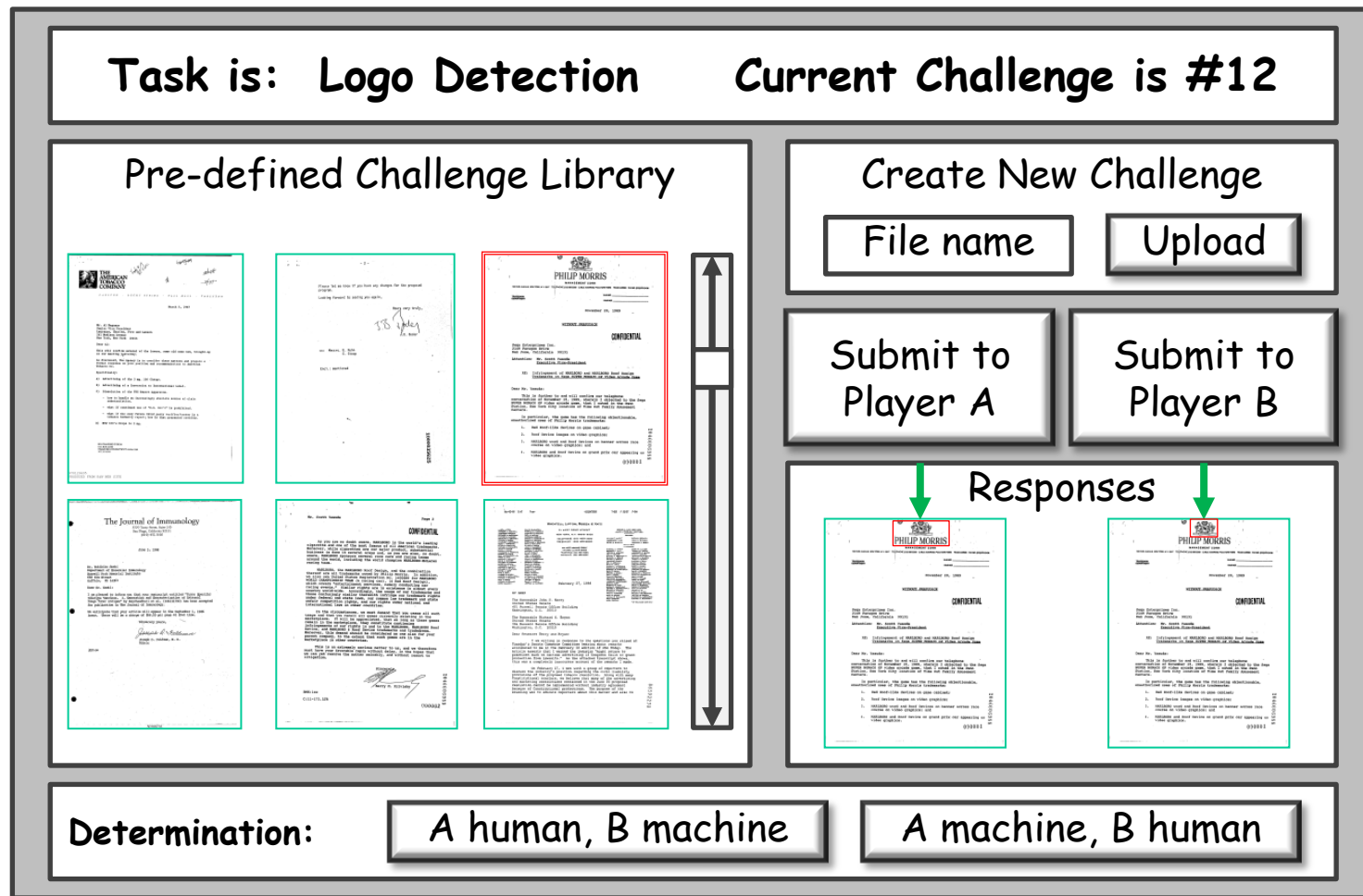
*Key Issues in Performance Evaluation for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 73*

# Other Considerations

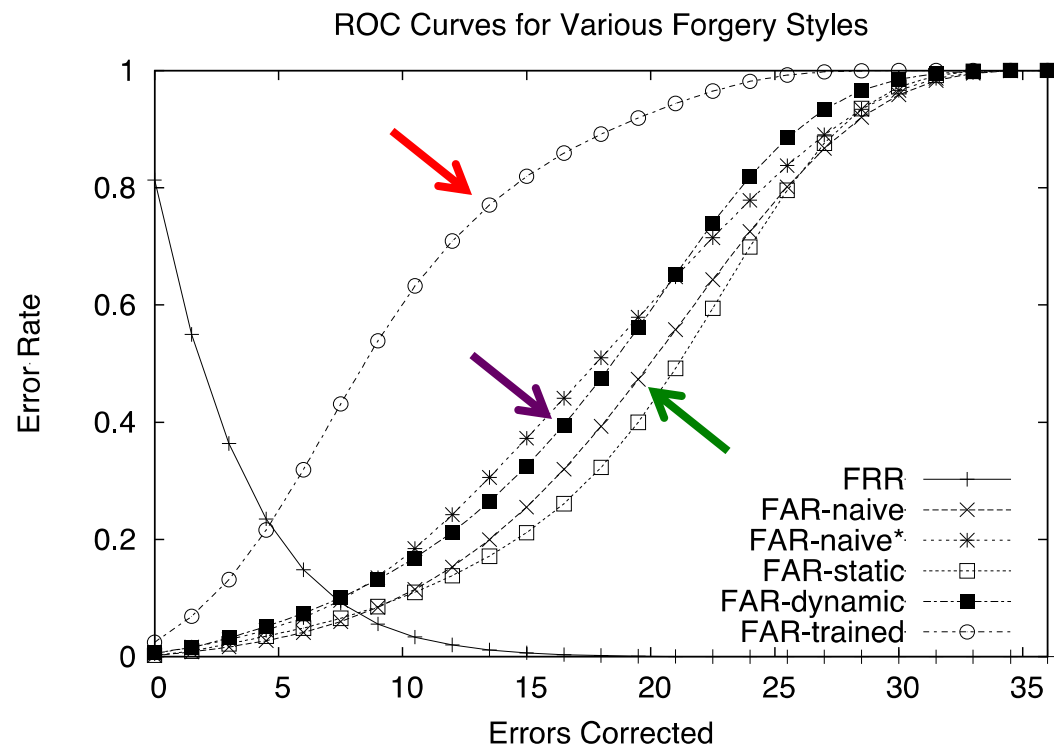Additional details to be addressed, some easy, some hard:

- Anyone should be permitted to volunteer at any point in time to serve as the judge or the human player.

- The need to pair a judge with a human player can be addressed through crowdsourcing (e.g., using micro-payments to recruit subjects like Mechanical Turk).

- How can we eliminate out-of-scope querying / collusion?

- Which problems are appropriate to test this way? (Avoid tedious tasks where machines are "too good.")

- How can learning (by human, by machine) be included?

"Adapting the Turing Test for Declaring Document Analysis Problems Solved," D. Lopresti and G. Nagy, Proceedings of the Tenth IAPR International Workshop on Document Analysis Systems (DAS 2012), March 2012, Gold Coast, Australia, 5 pages.

*Key Issues in Performance Evaluation*
*for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 74*

# Attack Models for Biometrics

My interests in rigorous, real-world performance evaluation have included research on attack models for behavioral biometrics, including online handwriting.

Trained, talented forgers are far more effective than "naïve" forgers, who are even bested by an automated synthesis technique.

ROC Curves for Various Forgery Styles

Error Rate vs Errors Corrected

Legend:
- FRR
- FAR-naive
- FAR-naive*
- FAR-static
- FAR-dynamic
- FAR-trained

"Forgery Quality and Its Implications for Behavioral Biometric Security," L. Ballard, D. Lopresti, and F. Monrose, IEEE Transactions on Systems, Man, and Cybernetics Part B, vol. 37, no. 5, October 2007, pp. 1107-1118.

*Key Issues in Performance Evaluation for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 75*

# Concluding Observations

- Play close attention to performance evaluation – it's important and not as straightforward as it may seem.

- Simply following common practice is not always enough.

- In most cases, ultimate goal is to replicate human interpretation for a pattern recognition task of interest.

- Recent developments – including new and better classifier technologies as well as the era of "big data" have led to tremendous breakthroughs and useful systems – but this doesn't diminish importance of performance evaluation.

- My thinking developed through collaborations with my students and colleagues, including Prof. George Nagy.

*Key Issues in Performance Evaluation for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 76*

# Thank you!

# ¡Gracias!

Key Issues in Performance Evaluation
for Document Analysis Systems

MCPR 2015
June 24-27
Mexico City, Mexico

Lopresti
Slide 77

# A Few Words About My University

Key Issues in Performance Evaluation
for Document Analysis Systems

MCPR 2015
June 24-27
Mexico City, Mexico

Lopresti
Slide 78

# Lehigh University



Packard Lab:  Home of
Computer Science & Engineering

- Private research university (1865)
- Four colleges:  Engineering, Arts & Sciences, Business, Education
- 441 full-time faculty members
- 4,577 undergrads, 2,064 grad students
- Three campuses, over 1,600 acres (side and top of mountain, heavily wooded)
- Located about 1.5 hours from NYC and Philadelphia, 3 hours from Washington
- Ranked in top 15% of U.S. national universities
- Ranked in top 20% of U.S. PhD-granting schools for engineering

MCPR 2015
June 24-27
Mexico City, Mexico

# Lehigh University



Lehigh University

120 km

New York

80 km

Philadelphia

Please pay me a visit if you're in the area!

*Key Issues in Performance Evaluation*
*for Document Analysis Systems*

MCPR 2015
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 80*

# CSE Department



Currently 15 faculty ...
BS, MS, PhD in CS and CompE:

- Bioinformatics
- Biomedical Image Analysis
- Data Mining
- Database Systems
- Document Analysis
- Intelligent Agents

- Mobile Robotics
- Networking
- Parallel Processing
- Programming Languages
- Computer Security
- Semantic Web
- Social Networking
- Web Search / Systems

*Key Issues in Performance Evaluation*
*for Document Analysis Systems*

**MCPR 2015**
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 81*

# Data X Strategic Initiative



http://www.lehigh.edu/datax

*Key Issues in Performance Evaluation*
*for Document Analysis Systems*

**MCPR 2015**
June 24-27
Mexico City, Mexico

*Lopresti*
*Slide 82*