

A New Paradigm for Experimental Pattern Recognition Research

(including Document Analysis and Exploitation*)

Daniel Lopresti

Computer Science & Engineering
Lehigh University
Bethlehem, PA, USA

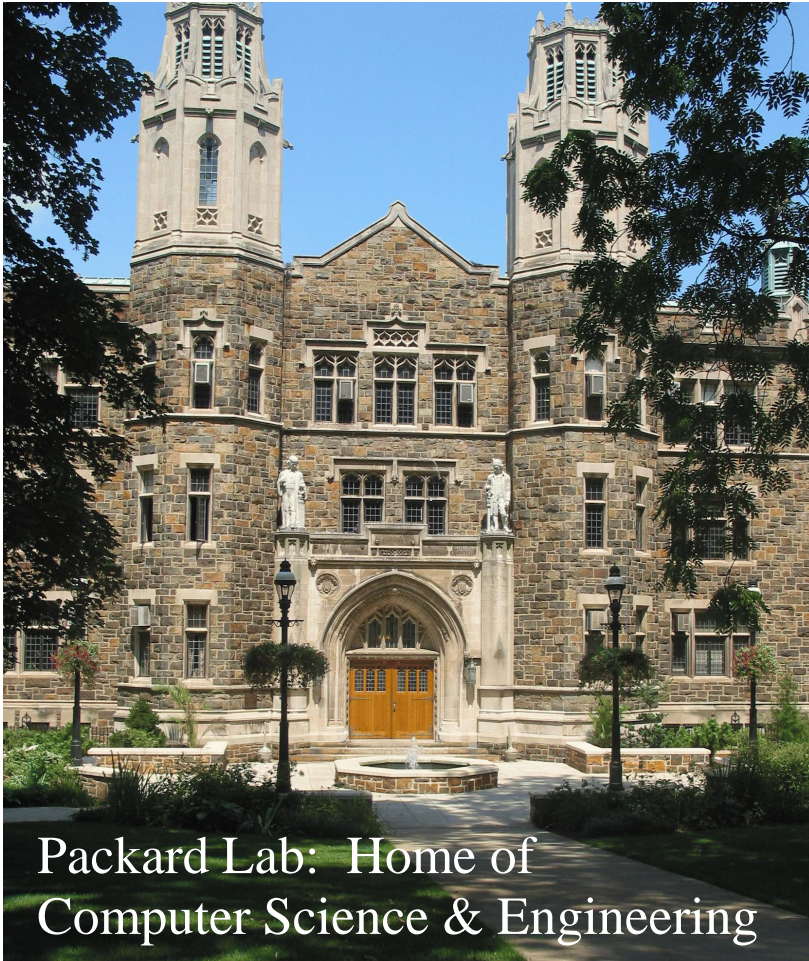
Bart Lamiroy

Nancy Université
Campus Scientifique, BP 239, 54506
Vandoeuvre Cedex, France

* Supported by a Congressional appropriation administered through DARPA IPTO via Raytheon BBN Technologies.



Lehigh University



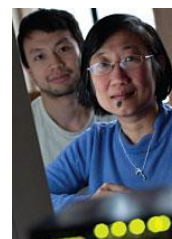
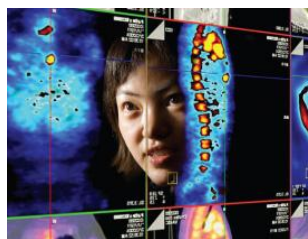
- Private research university (1865)
- Four colleges: Engineering, Arts & Sciences, Business, Education
- Faculty = 441 full-time
- Students = 2,064 grad, 4,577 undergrad
- Three campuses spread over 1,600 acres (mountain side, wooded)
- Located about 1.5 hours from NYC and Philadelphia, 3 hours from Washington
- Engineering College ranked in top 20% of Ph.D.-granting schools in U.S.
- University ranked in top 15% of U.S. national universities



Lehigh University



CSE Department



17 faculty, 60 PhD, 20 MSE ...

- Bioinformatics
- Biomedical Image Analysis
- Computer Game Design
- Data Mining
- Database Systems
- Document Analysis
- Intelligent Agents
- Mobile Robotics
- Networking
- Parallel Processing
- Programming Languages
- Computer Security
- Semantic Web
- Social Networking
- Web Search and Web Systems



Germ of an Idea (1)



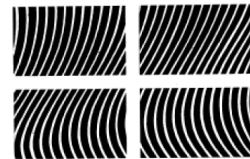
George Nagy graduated in Physics (fencing and solving Euler's Second problem) awarded the PhD at Columbia University. He worked with Rosenblatt building the Tobe network for speech recognition. He claims credit for IBM's reverse sabbatical at the University of Toronto. Department of Computer Science, Troy, NY. In 1985 he has been Professor of Computer Science at Troy, NY.

Nagy's credits in document analysis include: "self-corrective" character recognition with Henry Baird, character recognition with Henry Baird, character recognition with Henry Baird, character recognition with Henry Baird.

Wonderful keynote on datasets by George Nagy at DAS 2010 ...

Data Sets advertised in IEEE Computer

January 1972 (6 data sets)



PATTERN RECOGNITION DATA BASES AVAILABLE

In order to encourage research in the field of pattern recognition, the IEEE Computer Society's Technical Committee on Pattern Recognition has begun collecting data bases from a variety of sources. These data bases, including substantial back-up documentation, may be ordered by using the form at the bottom of the page.

1.1.1 Machine Imprinted Alphanumeric Characters - Dr. H. F. Row, Cornell Aeronautical Laboratory, Ithaca, N.Y., USA
An alphanumeric character data base of 100,000 samples of 48 character classes. (Also known as the CAL-U.S. Postal Service Alphanumeric Character Data Base). Threshold binary images of segmented, centered, mixed-font, machine-imprinted characters. Resolution is 24 x 24. Magnetic tape, 9 track, 2 reels, 1600 BPI.
Price: \$112.50 (\$68.75 with furnished tapes)
Member's discount price: \$90, \$55, with furnished tapes

4-COMPUTER/JANUARY/FEBRUARY 1972

Discounts off the data base list prices are available to IEEE members and members of the American Federation of Information Processing Societies' constituent societies.
When ordering, you may elect to send us your own blank tapes, if you do, be sure they are in good condition and have no other data recorded on them.

1.1.2 Handprinted Numeric Characters - Dr. A. L. Kriegl, Honeywell Information Systems, Data Systems Division
The data base consists of 50 samples of each numeric character generated by 9 different authors. Simple printing rules were specified but not always followed. The samples were selected from those contributed. The images are binary with a resolution of 25 x 21. Punched cards.
Price: \$37.50
Member's discount price: \$30.

PATTERN RECOGNITION DATA BASES AVAILABLE



but not always followed. The samples were selected from those contributed. The images are binary with a resolution of 25 x 21. Punched cards.
Price: \$41.25
Member's discount price: \$33.

1.1.2 Handprinted FORTAN Alphanumeric Characters - Dr. John H. Nelson, Stanford Research Institute
The data base consists of two parts, with each part to itself. The first part contains 11 alphabets of 46 characters, corresponding to the numeric character set of the mobile KODAK image, transmitted to each of 1000 users in a total of 3 x 64 x 49 x 1.25 pages.
The second part has 2,099 characters printed by a single author. There are 300 characters made up of 30 alphabets of 46 characters each, the remaining 2,079 characters are taken from fragments of actual coding teletype. The images are binary with a 24 x 24 resolution. Magnetic tape, 7 track, 2 reels, 1600 BPI.
Price: \$140.75 (\$74.75 with furnished tapes)

alphanumeric characters. The images are binary with a resolution of 12 x 12. Punched cards.
Price: \$54.
Member's discount price: \$44.

1.1.4 Handprinted Numeric Characters - Brooks Gordon, Telex-Behrens Systems Co., Ltd/Toshiba Research and Development Center
The data base consists of 10,000 hand written numeric characters collected from 5000 users in 1000 different locations. The data base is contained on two magnetic tapes with each tape having 1,000 characters. Images are binary with a resolution of 36 x 36. A single pattern consists of 50 words. Magnetic tape, 7 track, 2 reels, 56,000, 5.75 per character.
Price: \$182.25 (\$67.75 with furnished tapes)
Member's discount price: \$43, \$34, with furnished tapes

1.1.5 Courier Style - Dr. L. D. Hanson, Bell Telephone Laboratories

Pattern Recognition DATA BASES

1.1.1 Machine Imprinted Alphanumeric Characters - Dr. H. F. Row, Cornell Aeronautical Laboratory, Ithaca, N.Y., USA
An alphanumeric character data base of 100,000 samples of 48 character classes. (Also known as the CAL-U.S. Postal Service Alphanumeric Character Data Base). Threshold binary images of segmented, centered, mixed-font, machine-imprinted characters. Resolution is 24 x 24. Magnetic tape, 9 track, 2 reels, 1600 BPI.
Price: \$112.50 (\$68.75 with furnished tapes)
Member's discount price: \$90, \$55, with furnished tapes

1.1.1.1 Machine Imprinted Alphanumeric Characters - Dr. H. F. Row, Cornell Aeronautical Laboratory, Ithaca, N.Y., USA
An alphanumeric character data base of 100,000 samples of 48 character classes. (Also known as the CAL-U.S. Postal Service Alphanumeric Character Data Base). Threshold binary images of segmented, centered, mixed-font, machine-imprinted characters. Resolution is 24 x 24. Magnetic tape, 9 track, 2 reels, 1600 BPI.
Price: \$112.50 (\$68.75 with furnished tapes)
Member's discount price: \$90, \$55, with furnished tapes



To order: Use the coupon order form at the bottom of the page. Please specify to open with each tape having 1,000 characters. Images are binary with a resolution of 36 x 36. A single pattern consists of 50 words.

AS 2010 (GN)

April 1976 (10 data sets)

* Slides available on the DAS 2010 website.



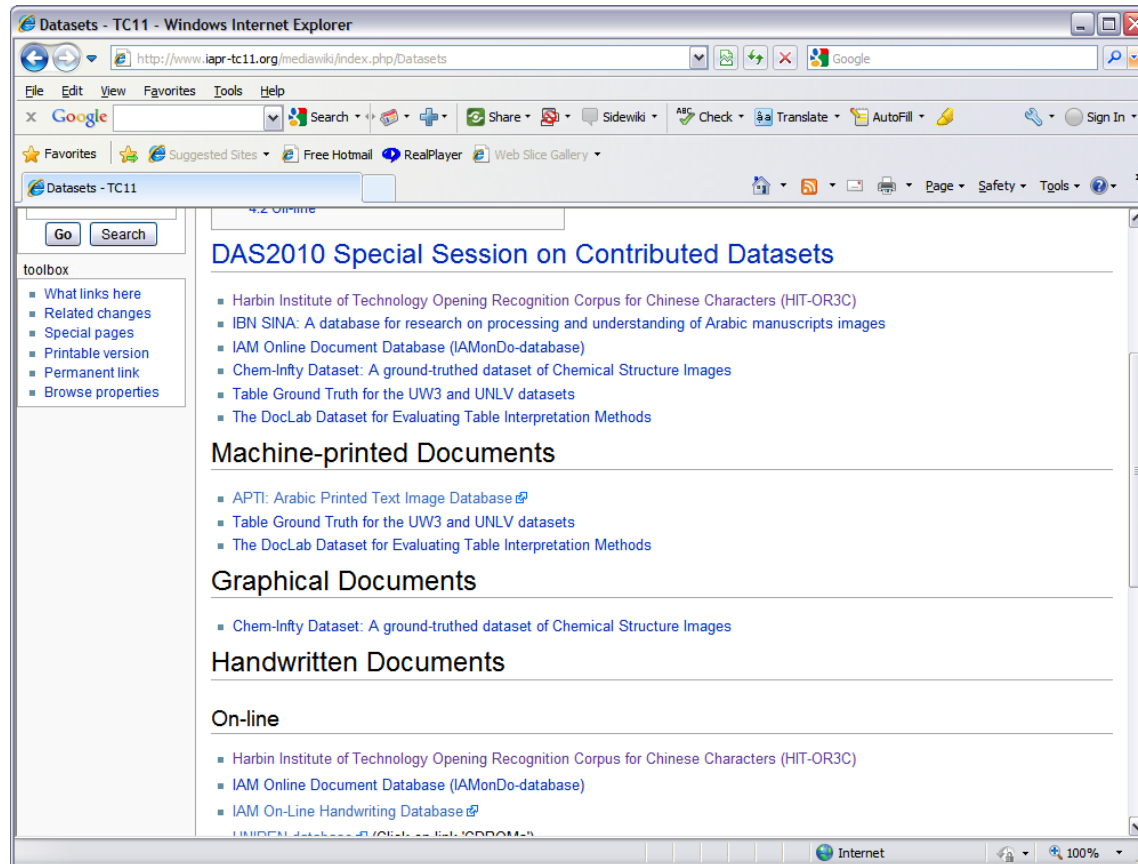
A New Paradigm for Experimental Pattern Recognition Research



Lopresti and Lamiroy September 2010 • Slide 5

Germ of an Idea (2)

IAPR TC-11 (“Reading Systems”) website:



<http://www.iapr-tc11.org>



Status Quo

Despite tremendous advances in computer and communications technologies, the way we conduct research and disseminate results is essentially unchanged over the past 50 years.*

Identify problem → conceive solution → find data →
test method → publish results (→ field system)

Modest steps forward include more stringent peer review processes (comparisons to prior work), attempts to create and share common datasets, and “competitions” at international conferences.

* See the excellent keynote talk delivered by George Nagy at DAS 2010:

http://cubs.buffalo.edu/DAS2010/GN_testing_DAS_10.pdf



Broad Aims for Research

Research goals for work in machine perception*, broadly-stated:

- Develop algorithms that are robust and approach human levels of performance for specific tasks of interest.
- Invent new methods that are better than known techniques.
- Generate experimental results that are well-documented, understandable in context, and reproducible by others.
- Build on past knowledge to yield new knowledge which moves us toward solutions for problems of vital importance.

* *Machine perception* refers to any pattern recognition task that attempts to mimic human behavior, including computer vision, document analysis, biomedical image processing, speech recognition, natural language understanding, etc.



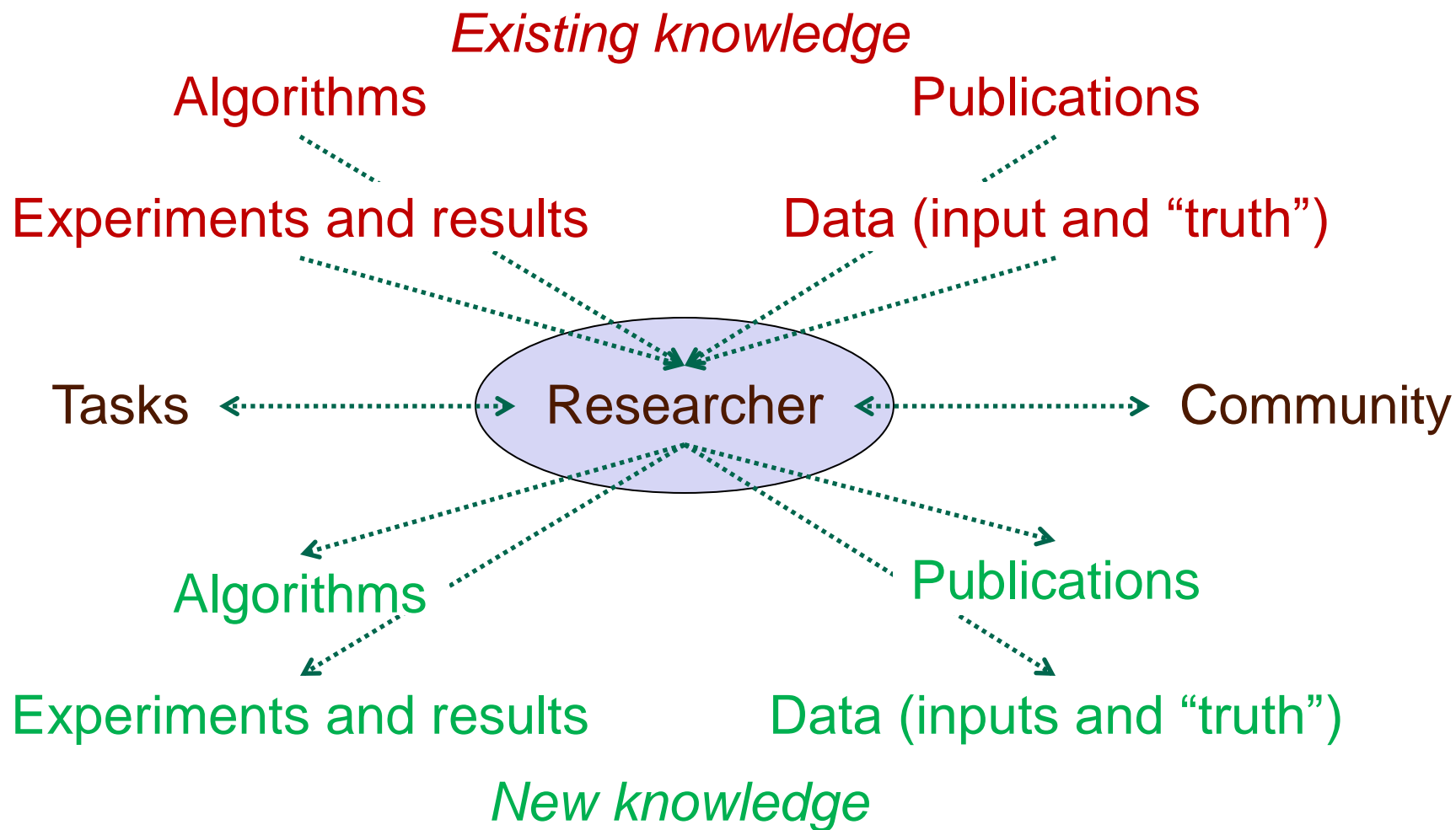
Moving Forward

A fundamentally new paradigm for experimental research:

- Not a minor “tweak” – an enormous, ground-breaking change.
- Enabled by infrastructure that largely exists – most pieces are already present, they just need to be put into place.
- Could dramatically accelerate the rate of scientific discovery.
- Power derived from community participation, but vision is so compelling that some version of this future is inevitable.
- Critical mass of early adopters (thought-leaders) will play a key role in guiding the vision to fruition. Despite open-access model, significant proprietary benefits are also possible.
- Large potential for “meta” research (developing the paradigm).



Current Paradigm



Problems with Status Quo (1)

Stated goal:

- Develop algorithms that are robust and approach human levels of performance for specific tasks of interest.

Observations:

- We want algorithms to be general, but too often they are tested on small, overused datasets far removed from the real world.
- Nearly all experimental results reported in the literature are biased by algorithm developer's intimate knowledge of the data.
- Current practices lack convincing evidence of generality.
- What does “human levels of performance” mean? Even experts can disagree on all but the most trivial of cases.



Problems with Status Quo (2)

Stated goal:

- Invent new methods that are better than known techniques.

Observations:

- How do we know when we have succeeded?
- Need to compare against previously published results creates over-reliance on standard datasets (which is counter-productive).
- Attempts to re-implement a published algorithm are problematic (incomplete descriptions, inherent conflict of interest).
- Competitions can be useful, but are infrequent.
- Most papers do not even bother to make such comparisons.



Problems with Status Quo (3)

Stated goal:

- Generate experimental results that are well-documented, understandable in context, and reproducible by others.

Observations:

- Are published descriptions sufficient to reproduce experiments?
How often is this even attempted?
- Explicit and/or implicit bias in selecting and using data (e.g., discarding hard cases) makes context difficult to recover.
- “Publish or perish” mindset leads to overstated claims, poor understanding of generalizability of results.



Problems with Status Quo (4)

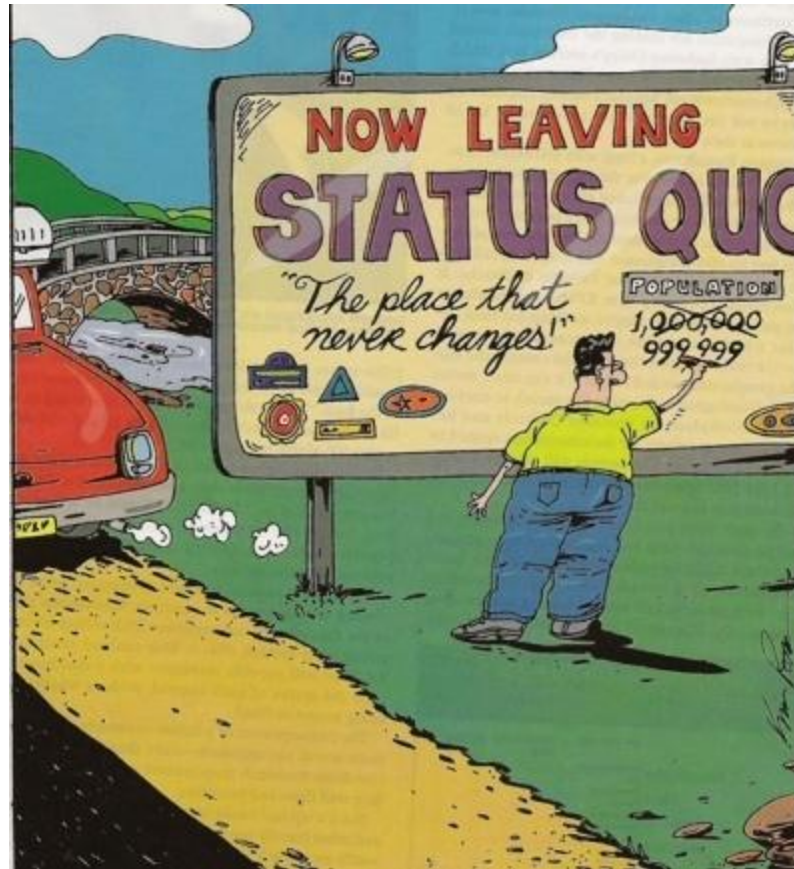
Stated goal:

- Build on past knowledge to yield new knowledge which advances the field toward solving problems of vital importance.

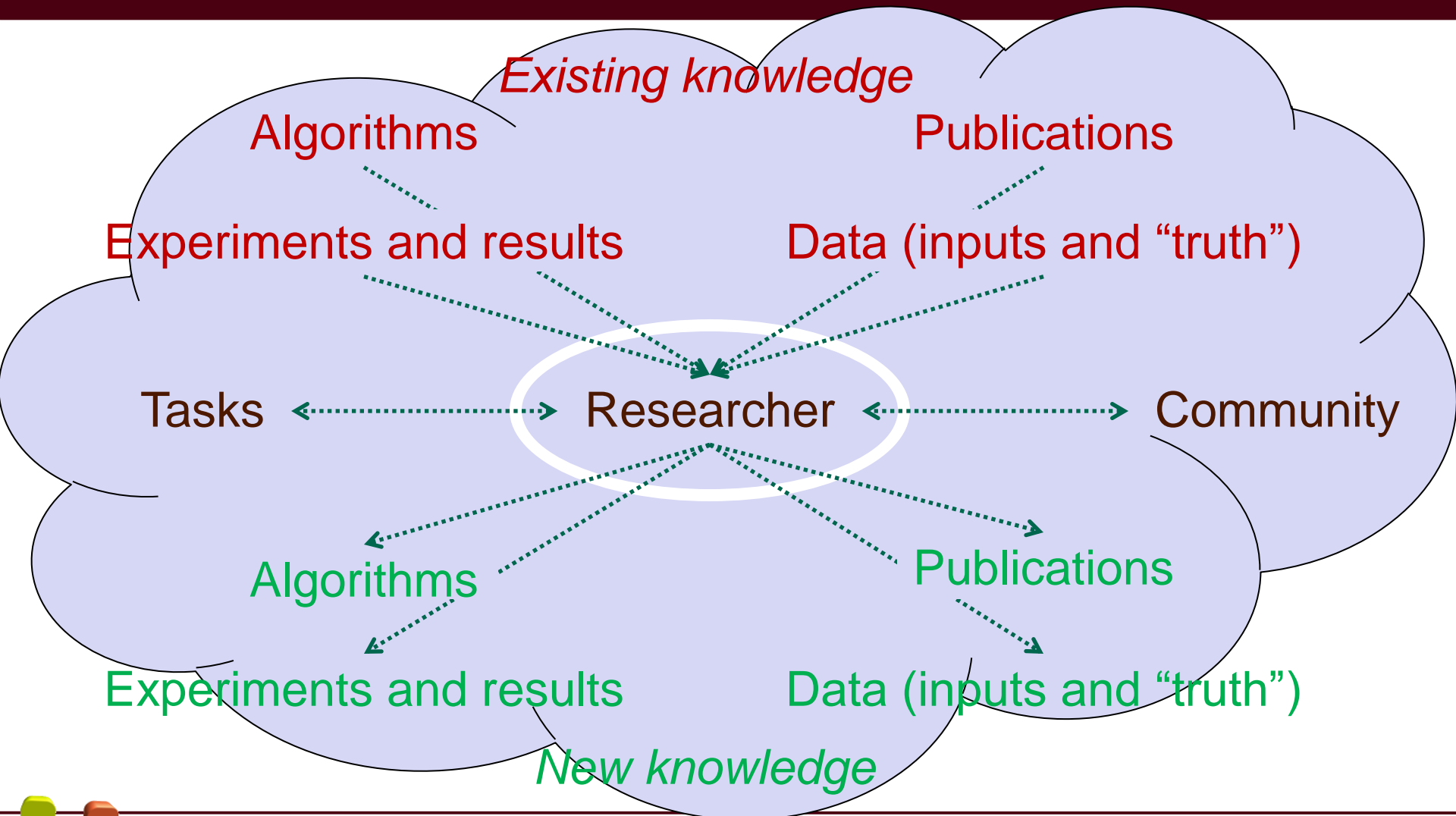
Observations:

- Effort is wasted developing methods that do not improve on existing techniques, or are ill-suited for the task at hand.
- Much time spent “reinventing the wheel.”
- Like trying to build a pyramid out of shifting sand without first forming it into blocks.
- Impossible to fix without a paradigm shift across the community.





Vision: A Communal Infrastructure

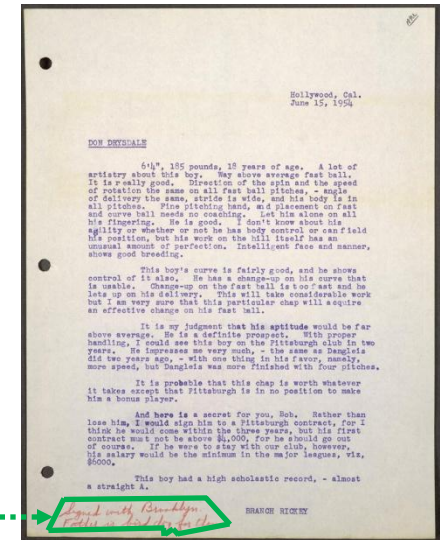


Looking Ahead: A Scenario



Sometime in 2020, Jane, a young researcher just getting started in the field turns her attention to a specific task: given a page image, identify regions that contain handwritten notations.*

- Jane wants a fully general method that can accept any scanned page as input.
- She plans to develop an algorithm that will place a bounding polygon as tight as possible around each handwritten region.
- Each polygon should enclose a logical “unit.”



* If you prefer, replace this task by any one in machine perception you like.



Truth vs. Interpretation

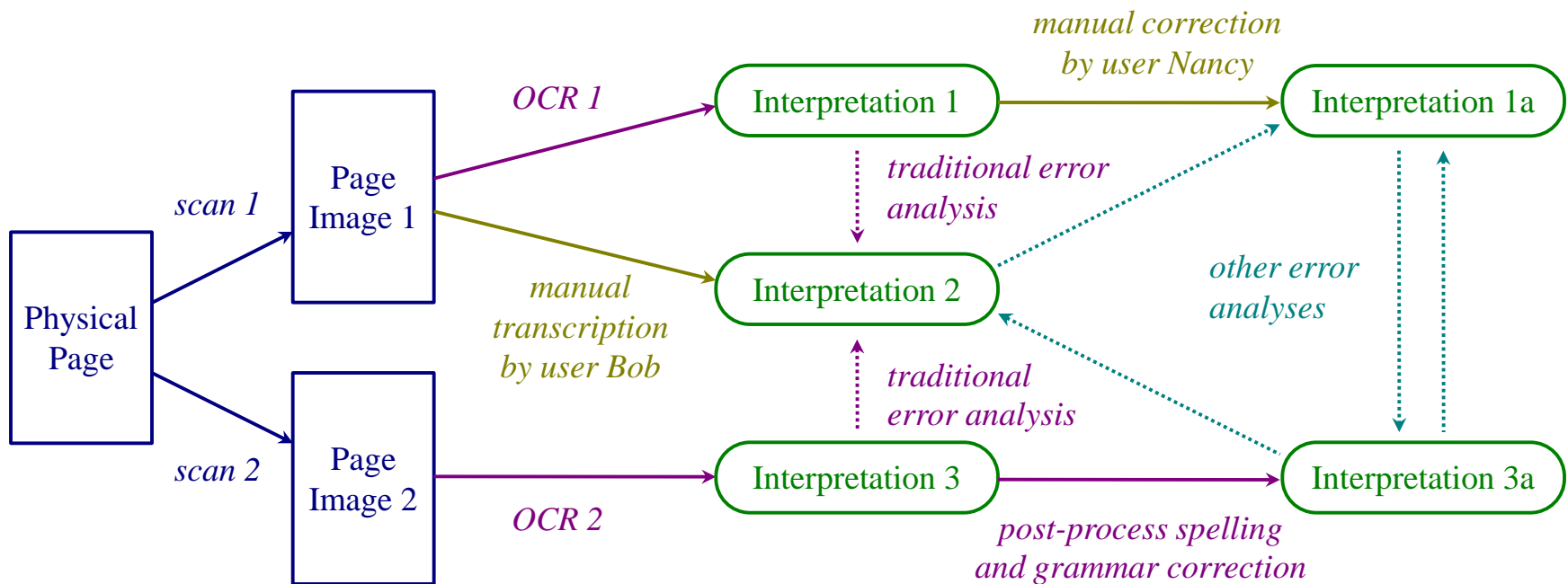
What should Jane be aiming for?

- In this version of the future, there is no such thing as “ground truth” – that term is no longer used.
- Rather, we talk about the *intent* of the author (which can be hard to determine, although sometimes we have access to it) ...
- ... or the *interpretation* arrived at by a reader of the document (which could be a human or an algorithm).
- There are no “right” or “wrong” answers – interpretations may naturally differ – although for some applications, we expect that users who are fluent in the language and the domain of the document will agree nearly all of the time.



Alternate Interpretations

Entities on a page are located and interpreted by humans interacting with the system, as well as by algorithms which have been invoked to process the page. This raises the possibility of alternate interpretations and relativistic error analyses:

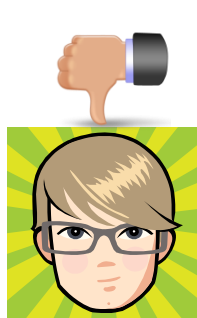


The Role of Reputation

With multiple competing interpretations, how does Jane proceed to develop a new method that mimics a careful, fluent human expert?

- Some people are more careful, fluent, and/or expert than others.
- The concept of online reputation as originally derived for use in social networking (Web 2.0) determines whose interpretations Jane will trust when developing her new method.
- This use of reputation is one key feature of the new paradigm.

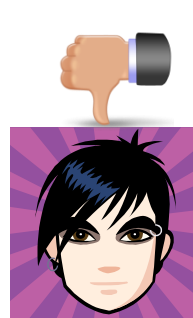
Users who have contributed interpretations



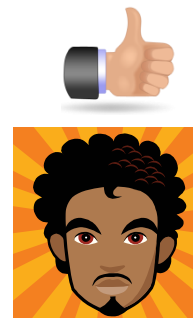
Bob



Alice



Hideo



Chuck



Kim



Kinds of Reputation

We have always had a notion of reputation, even in the early days. What is different in 2020 is this has been formalized to support experimental research.



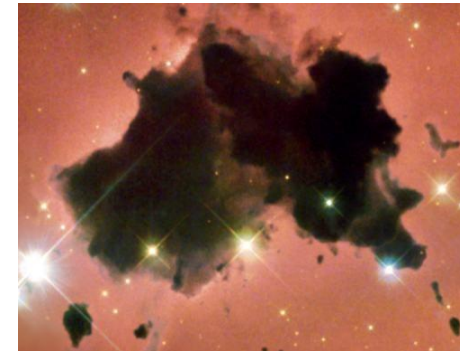
Various entities and possible bases for their reputations:

- Algorithm (how well does it address a stated task?)
- Implementation (is it buggy, or robust to a variety of inputs?)
- Dataset (is it representative for a particular problem?)
- Interpretation (is it trustworthy?)
- Publication (is it highly cited?)
- Researcher (amalgamation of contributions to above)



The DARE Server

After Jane has defined her task, she turns to a web server called DARE (for “Document Analysis Research Engine”), to request a set of sample pages for use in developing her new algorithm.



- DARE server lives in the “cloud” – it is not a single machine.
- A community-managed resource which has become a de facto standard (just as UW-1 dataset was once considered a standard).
- Collects datasets, interpretations, algorithms, citations, etc.
- Manages reputations and provides a query interface and API.
- Not just a storage system – DARE supports our new paradigm for experimental document analysis research.



Building Datasets via Queries

Jane queries the server to give her 1,000 random pages from the various collections it knows about. Through the query interface, she specifies:



- Each page should be 300 dpi color TIF. Jane may allow conversion from other formats (via built-in utilities), or she may prohibit this.
- Each page should be “predominately” printed material: text, line art, photographs, etc. This might mean $\geq 60\%$ of the page area.
- A reasonable number of pages should contain at least one annotation. For Jane’s application, this might be 10% to 25%.
- The handwritten annotations should be delimited in a way consistent with the intended output from Jane's algorithm.



The File Format Issue

In the past, there were attempts to unify the research community around common file formats. Since there was no compelling benefit to the individual researcher, these were routinely ignored.

- The status of the DARE web server as the de facto standard has enticed most researchers to use compatible file formats.
- DARE is the preferred platform to support their research.
- There is no coercion – it is just easier for users to go along since so much data is now delivered from the server and no longer lives locally on their own machines anyway.
- Some find this analogous to the earlier rise of other standard file formats such as TIF, JPG, PDF, PS, RTF, etc.



Logging and Transparency

In processing Jane's request, the DARE server returns the set of 1,000 random pages along with their associated interpretations.

- It also makes a permanent record of her query and provides a URL that will return exactly the same set of pages each time it is run.
- Any user with the URL can access the dataset and see Jane's parameter settings.
- All accesses are logged.

<http://www.dareserver.org/query?ID=124881371744>

Dataset created by Jane Simmons on July 26, 2020.

Parameter settings:

- 1,000 random page images
- 300 dpi color TIF format (no conversions)
- At least 60% foreground (printed) material
- 10% to 25% containing at least one annotation
- XMLpoly interpretation format
- Highest available reputation for data and interpreter

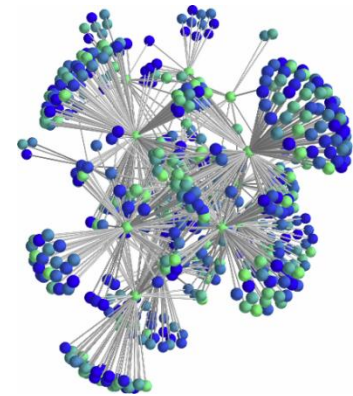
Click [here](#) to download the dataset.

Click [here](#) to access the log for this dataset.



Resource Discovery

Based on Jane's browsing behavior on the server, and her eventual query to create a dataset, DARE is smart enough to glean what she may be doing. Other people have made similar queries in the past.

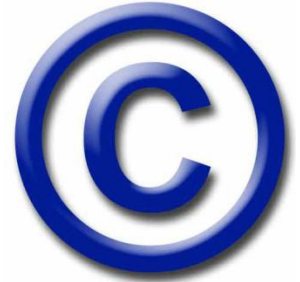


- DARE suggests a number of previously published papers addressing similar problems. Jane already knows about some of these, but others are new to her. She notes this for reference.
- The server also reports that several researchers have registered algorithms that Jane may choose to test her own against later. (For security and reliability, the DARE server runs contributed code using a virtual machine architecture.)



The Copyright Question

In the field's early days, the sharing of datasets was hampered by copyright issues. A group could create data for its own use, but was prevented from sharing it with others. By 2020, this is no longer a problem:



- Vast quantities of image data are available on the WWW. Digital libraries, both commercial and public, contain billions of pages.
- OCR and manual transcriptions have made this data searchable via simple text queries, but more sophisticated methods are needed to deal with “information overload” (hence Jane’s work).
- Some digital libraries require a subscription. The DARE server knows which libraries Jane has access to and manages the authentication process for her seamlessly.



Other Kinds of Datasets

The DARE server is not limited to delivering randomly-generated datasets. It supports other contributions from the community including:



- Legacy datasets like UW-1 referenced by their own URL's:
<http://www.dareserver.org/dataset?ID=UW-1>
- Components available for synthesizing page images using pre-defined models for layout, typography, and noise effects.
- A “black box” testing protocol for data that is stored securely on the DARE server but that cannot be re-distributed.



The Creative Process is Unchanged

With dataset in hand, Jane starts to work on her algorithm. This step is no different from what researchers do today.

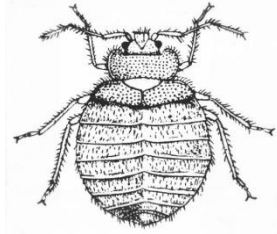


- Jane may study the pages in the dataset provided by the server to better understand the cases her method should handle.
- If her algorithm employs machine learning, she could use some of the pages as a training set and others as a “test” set.
- Since it is understood that the design of her algorithm may be biased by knowledge of the specific pages in the dataset, and since her goal is a more general method, Jane will only use this dataset for development purposes and never when it comes time to publish her results.



Buggy Data

The DARE server manages large quantities of data contributed by the community with no direct oversight. Hence, it is understood that bugs will arise occasionally.



- While working with the data, Jane notices a few problems. One of the page images was delivered to her upside down (rotated by 180 degrees). In another case, the TIF file for a page was unreadable, at least by the version of the library she is using.
- Being a responsible member of the research community (and wanting her online reputation to reflect this fact), Jane logs onto the DARE server and, with a few mouse clicks, reports both of these problems – it just takes a minute.
- Jane's reports will be checked by other community members.



Differing Interpretations

In a few other cases, Jane disagrees with the interpretation provided for the page in question.



- In her opinion, on one page the bounding polygons are drawn improperly and, on another, an annotation has been missed.
- Rather than just make local changes to her own private copies of the data (as would have happened in the past), Jane records her interpretations on the DARE server. Her opinions are added to the existing ones, enriching the entire collection.
- To facilitate this, the server provides a wiki interface with editing and markup tools that run in any web browser.



Unbiased Testing

After Jane is done fine-tuning her algorithm, she prepares to write a paper about it. This will involve testing the claim that her method works for arbitrary inputs, not just for the data she has been using for the past six months. The DARE server gives her two options for performing random, unbiased testing of her method. These are:

- Option #1 (local): Jane can “wrap” her code in a driver provided by the DARE server that runs the algorithm on a series of page images, calling the server each time a new page is required.
- Option #2 (remote): Jane can choose to upload her code to the DARE server, thereby contributing it to the community. The server runs her algorithm on a series of previously unseen page images according to her specifications.



Unbiased Testing (Local)

Jane wraps her code in a driver provided by the DARE server ...

- The code continues to run on Jane's machine, with the server delivering a series of random page images it has not seen before but that satisfy certain properties she has specified in advance.
- The algorithm performs its computations and responds to the server within a few seconds – too fast for a human to intercede.
- Testing takes place over a period of time, with random delays between when the next page is requested and when it is delivered. It would be too tedious for a human to try to “game” the system.
- Hence, everyone can be certain that the results reported are in fact the unbiased output from Jane's program.



Unbiased Testing (Remote)

Instead, Jane can choose to upload her code to the DARE server ...

- In this case, the server will run her algorithm on a series of previously unseen test pages according to her specifications.
- Since the testing is under the complete control of the DARE server, there is no risk of someone trying to “game” the system.
- A virtual machine architecture is used for security and reliability. Malicious code is contained and cannot compromise the server.
- The DARE server will also maintain Jane’s code in the system and use it in future comparisons when other researchers test their own new algorithms on the same task.



Local vs. Remote Testing Models

The two schemes, local vs. remote, provide tremendous flexibility:

- The first option allows Jane to keep her code private, if she so wishes, since the code runs entirely on her own machine. Perhaps she is still developing her method and would rather not publicize it until she can improve it, or maybe she works for a company or a government agency, or she is planning to file a patent application.
- The second option, of course, allows Jane to contribute her algorithm to the community (and, again, raise her online reputation), just as anyone can add new data or interpretations to the DARE server.



Certified Testing Results

At the end of the evaluation, Jane is provided with:

- A set of summary results showing how well her algorithm matched human performance on the task.
- Another set of summary results showing how well her algorithm fared in comparison to other methods tested on the same pages.
- A *certificate* (i.e., a unique URL) that guarantees the integrity of the results and which can be cited in the paper she is writing. Anyone who enters the certificate into a web browser can see the results delivered directly from the (trusted) DARE server.
- Now there can be no doubt that what Jane reports in her paper is true and scientifically reproducible.



Certified Testing Results



<http://www.dareserver.org/query?ID=08272348481991>

Certified evaluation run by Jane Simmons on December 3, 2020.

Parameter settings:

- 2,000 random page images guaranteed previously unseen
- 300 dpi color TIF format (no conversions)
- At least 60% foreground (printed) material
- 10% to 25% containing at least one annotation
- XMLpoly interpretation format
- Highest available reputation for data and interpreter

Guaranteed unbiased testing

Performance relative to human interpretation:

- 92% recall / 89% precision
- No result returned for 2% of inputs (click [here](#) to access these pages).

Comparison to human expert

Performance of best previously registered algorithm on same pages:

- 87% recall / 82% precision
- +5% net improvement in recall / +7% net improvement in precision

Comparison to earlier method

Click [here](#) to download the dataset.

Click [here](#) to access the log for this dataset.



Promoting Robustness

It was perhaps a bit optimistic of Jane to believe that her method would handle all possible inputs and, in fact, she learns that her code crashed on 2% of the test pages.

- The server allows Jane to download these pages to see what is wrong (she failed to dimension a certain array to be big enough).
- If Jane is requesting a certificate, the DARE server will guarantee that her code never sees the same page twice.
- If she is not requesting a certificate, there is no restriction and the server will deliver the same page as often as she wishes.
- Past researchers could remove troublesome inputs from their test sets, but the DARE server now prohibits this.



Task-Based Evaluation

In the past, it was well known that simple summary measures like recognition accuracy and precision / recall may not correlate with performance in a target application. Nothing could be done about this, however, since the extra effort required to implement a full end-to-end system was too much work for a single researcher.

- Now, in 2020, Jane can choose to evaluate her algorithm as part of a pre-programmed pipeline that consumes the annotations it identifies, performs handwriting recognition, and tests the results in prototypical information extraction and retrieval applications.
- As before, relative improvement over existing methods is key.
- A variety of pipelines are fielded on the server for testing.



Other Notable Benefits

The new paradigm brings with it a range of significant benefits:

- Experimental results are now verifiable and reproducible.
- Beginning researchers no longer waste their time pursuing methods that are inferior to known techniques (since the DARE server will immediately inform you if another registered algorithm did a better job on the test set you were given).
- The natural (often innocent) tendency to bias an algorithm based on knowing the details of the test set has been eliminated.
- The overuse of small “standard” datasets so prevalent in the early days of the field is now no longer a problem.



Questions Never Before Answerable

Jane's queries are just one example of what DARE makes possible. A system such as this will allow users to ask questions like:

- What is the best known algorithm for a certain task? (This is the method that comes closest to matching human performance.)
- What is the domain of expertise for a particular algorithmic technique? (These are the categories of inputs for which the method performs significantly better than it does elsewhere.)
- What classes are most vexing to document image analysis? (These are the inputs for which known methods perform poorly.)
- What is the “hardest” open problem in a given field of inquiry? (This is the task for which the current gap between machine and human performance is the largest.)



Attempts to Defeat the System

The DARE server is not foolproof – it provides many features to encourage and support good science, but it cannot completely eliminate the possibility of abuse by a malicious individual.

- For example, in order to guarantee that a given page is not seen twice, it is necessary to register both the user and the algorithm.
- Just as with any other web service, it is possible to create multiple accounts in an attempt to circumvent this rule.
- However, due to the community nature of DARE, all records are open and visible to every user of the system, increasing the risk of being caught to the degree that legitimate researchers would never be willing to take that chance.



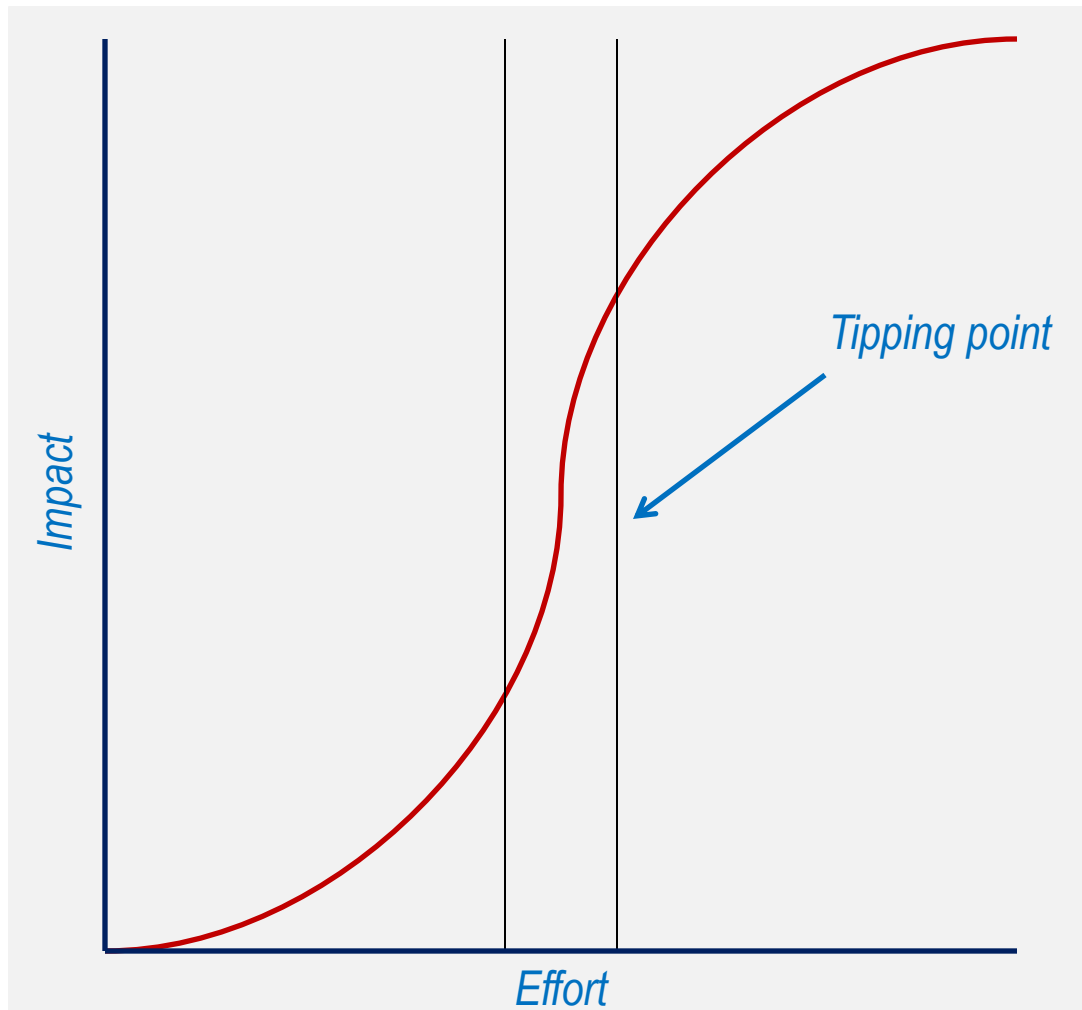
Key Pieces Already Present

Much of this should sound familiar because it is already happening elsewhere in other contexts (including Web 2.0):

- amazon.com: recommendations based on searches and past purchases, independent reviewers who in turn are rated by users.
- Netflix: movie suggestions based on ratings by >1M customers.
- Wikipedia: community authorship of a valuable online resource.
- YouTube and Flickr: user-provided and annotated multimedia.
- NLM model for interactive scientific publication.
- Work by NLM / Lehigh team on multi-observer segmentation.
- Recent research on data provenance (e.g., VisTrails system).
- IAPR TC-11 and TC-10 efforts to collect standard datasets.



Where is the Tipping Point?



\$64,000 question:
is required effort
modest or large?

We believe it should
be possible to
achieve key aspects
of what we have
described here with
a few year's work.



Possible Paths Forward

Currently ongoing:

- Lehigh DAE project.
- IAPR TC-11 and TC-10 dataset efforts.

Worth exploring:

- Grants, partnerships, evangelizing within the community.

From the international research community:

- Ideas?
- Interest??
- Involvement???

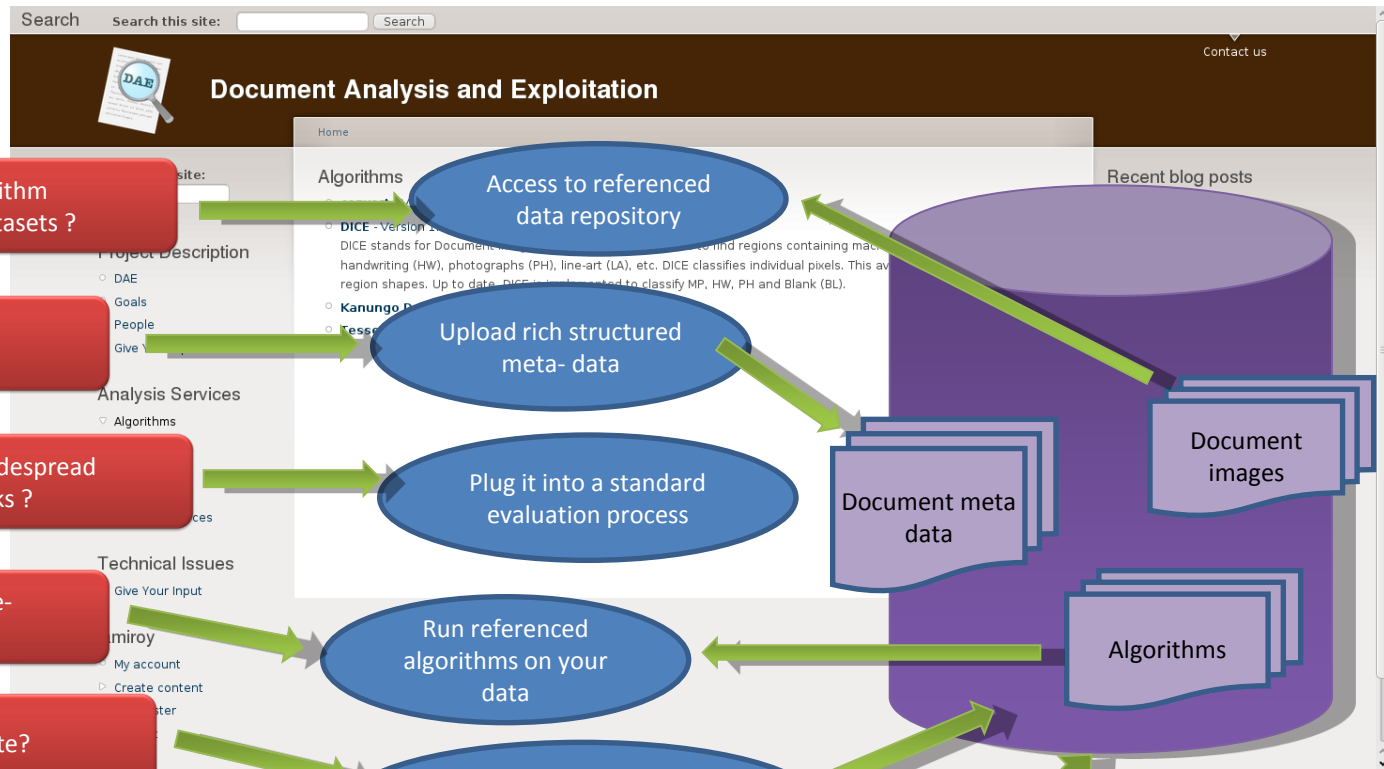


Why this Future is Inevitable

- Not everything discussed is under our control (e.g., copyright is an enormous legal issue), so why are we so certain about this?
- Small hurdles mean scenario will likely play out differently, but enough is doable that it is entirely feasible to make this happen.
- Focus on one key question: is the vision compelling and will it advance science? If the answer is “yes,” then it must happen.
- From a competitiveness standpoint, even though what we have described is a communal resource, some individuals / institutions / countries will be better positioned to exploit it at first than others.
- The tipping point may arrive sooner than we expect.



DAE Server (dae.cse.lehigh.edu/)



How does your algorithm behave on standard datasets ?

Add a new interpretation?

Are you compatible with widespread evaluation benchmarks ?

How do you compare to state-of-the-art algorithms for your problem ?

Do you want to compete?

Do you want to be peer-visible?

Disagree ?

Access to referenced data repository

Upload rich structured meta-data

Plug it into a standard evaluation process

Run referenced algorithms on your data

Upload your algorithm for execution by others

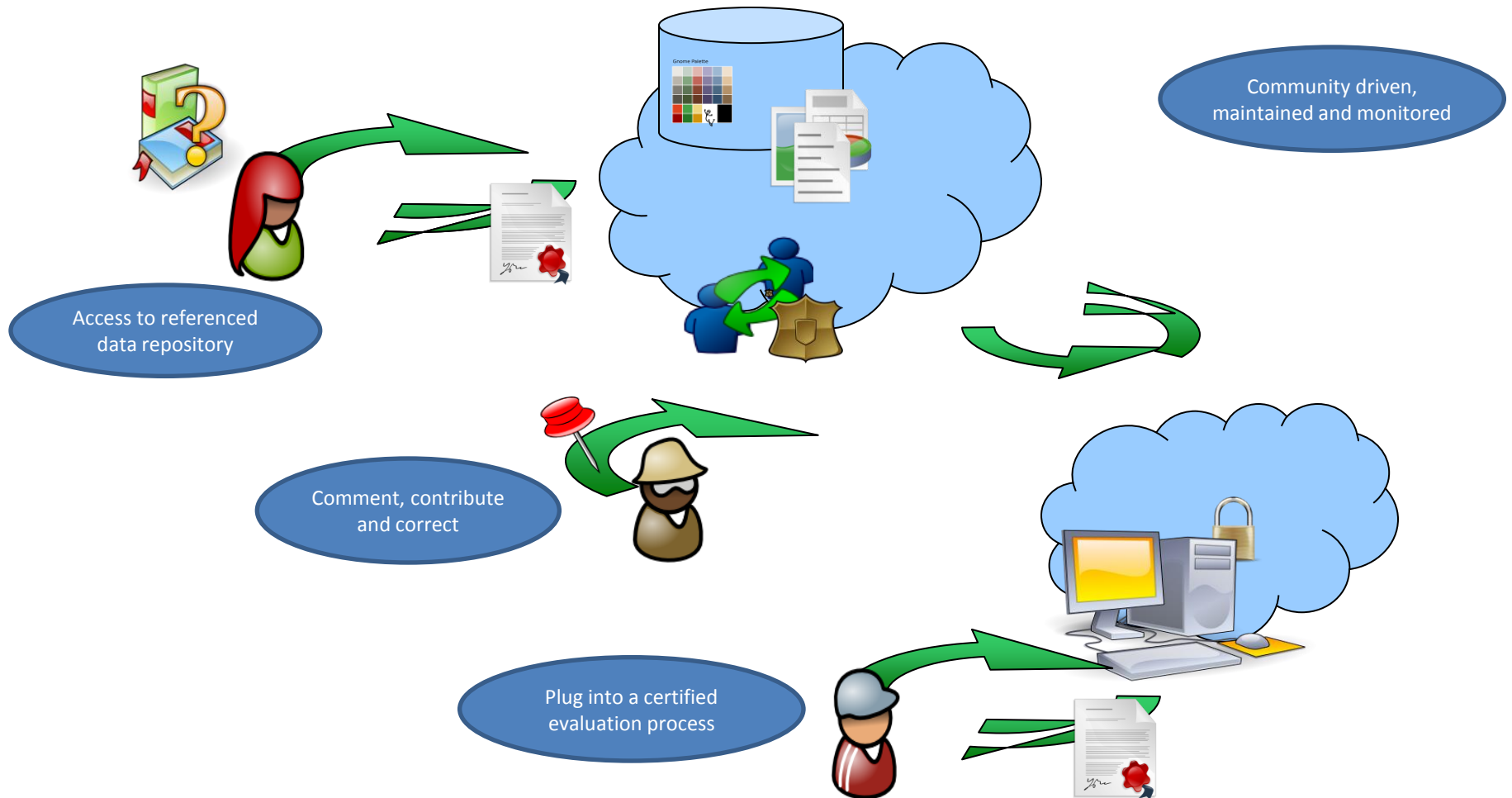
Comment, contribute and correct

Document meta data

Document images

Algorithms

DAE Server Technical Overview



Hardware & Software

Hardware:

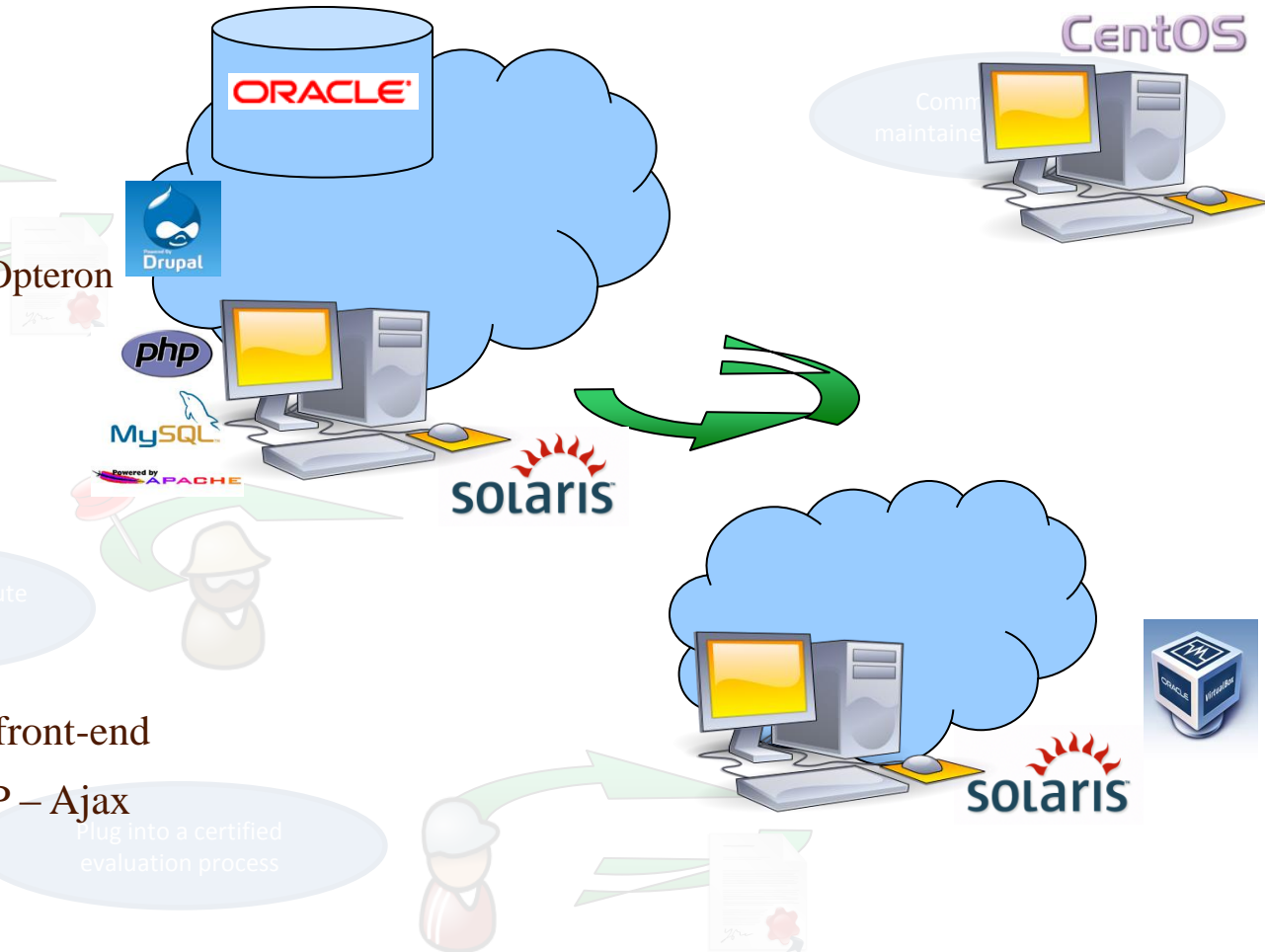
- Sun Fire X4540
 - 48TB raw storage,
 - dual six-core AMD Opteron

Architecture:

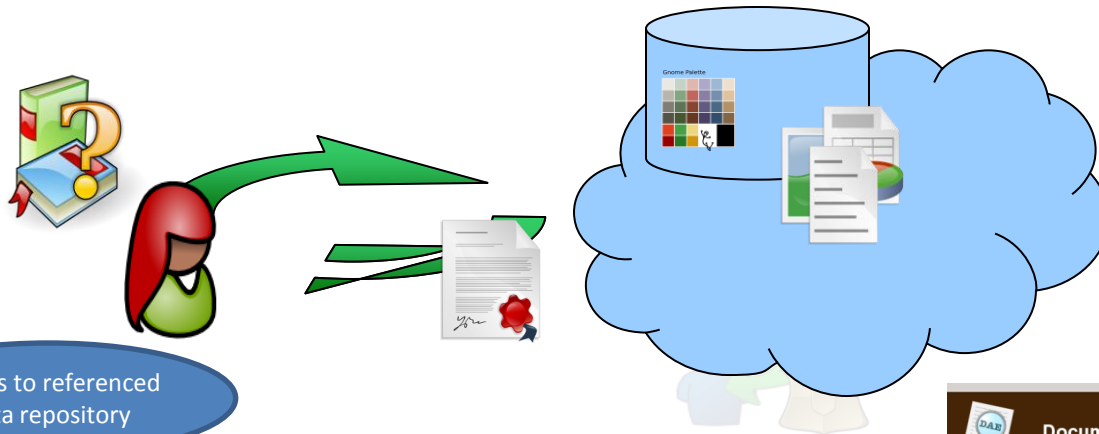
- Solaris Zones
- VirtualBox
- Beowulf HPC

Software:

- Drupal powered LAMP front-end
- Custom open-source PhP – Ajax
- Oracle back-end



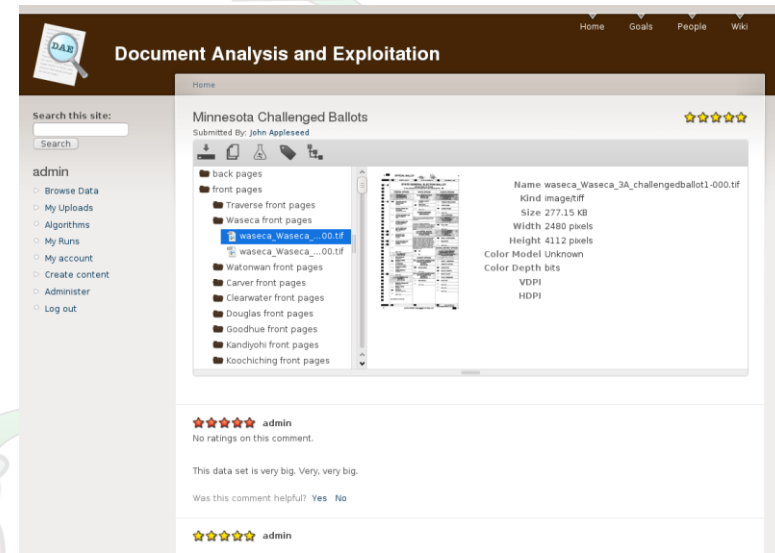
Intelligent Data Repository



Community driven,
maintained and monitored

Currently Implemented:

- Web-based browsing and downloading
- XML annotation format for uploading
- Full raw SQL interface
- Search by tag name



Intelligent Data Repository

```

<dataset>
  <id>3_G</id>
  <name>UNLV 3_G Sub-Dataset
    for testing</name>
  <description>
    For specific information about the contents
    consult the file(s) README.V3 on the appra
  </description>
</dataset>

<page_image>
  <id>ISRI-OCRTk/3/3_G/0148_479.3G</id>
  <path>ISRI-OCRTk/3/3_G/0148_479.3G</pat
  <vdpi>300</vdpi>
  <hdpi>300</hdpi>
  <skew></skew>
  <type_list>
    <type_list_item>
      <id>TIFF file, big-endian</id>
    </type_list_item>
    <type_list_item>
      <id>image/tiff</id>
    </type_list_item>
  </type_list>
  <contributor_list>
    <contributor_list_item>
      <id>isri-info@isri.unlv.edu</id>
    </CONTRIBUTOR_LIST_item>
  </contributor_list>
  <copyright>UNLV Copyright</copyright>
  <description>UNLV ISRI-OCRTk/3/3_G/0148
  <image_dataset_list>
    <image_dataset_list_item>

```

Download Button

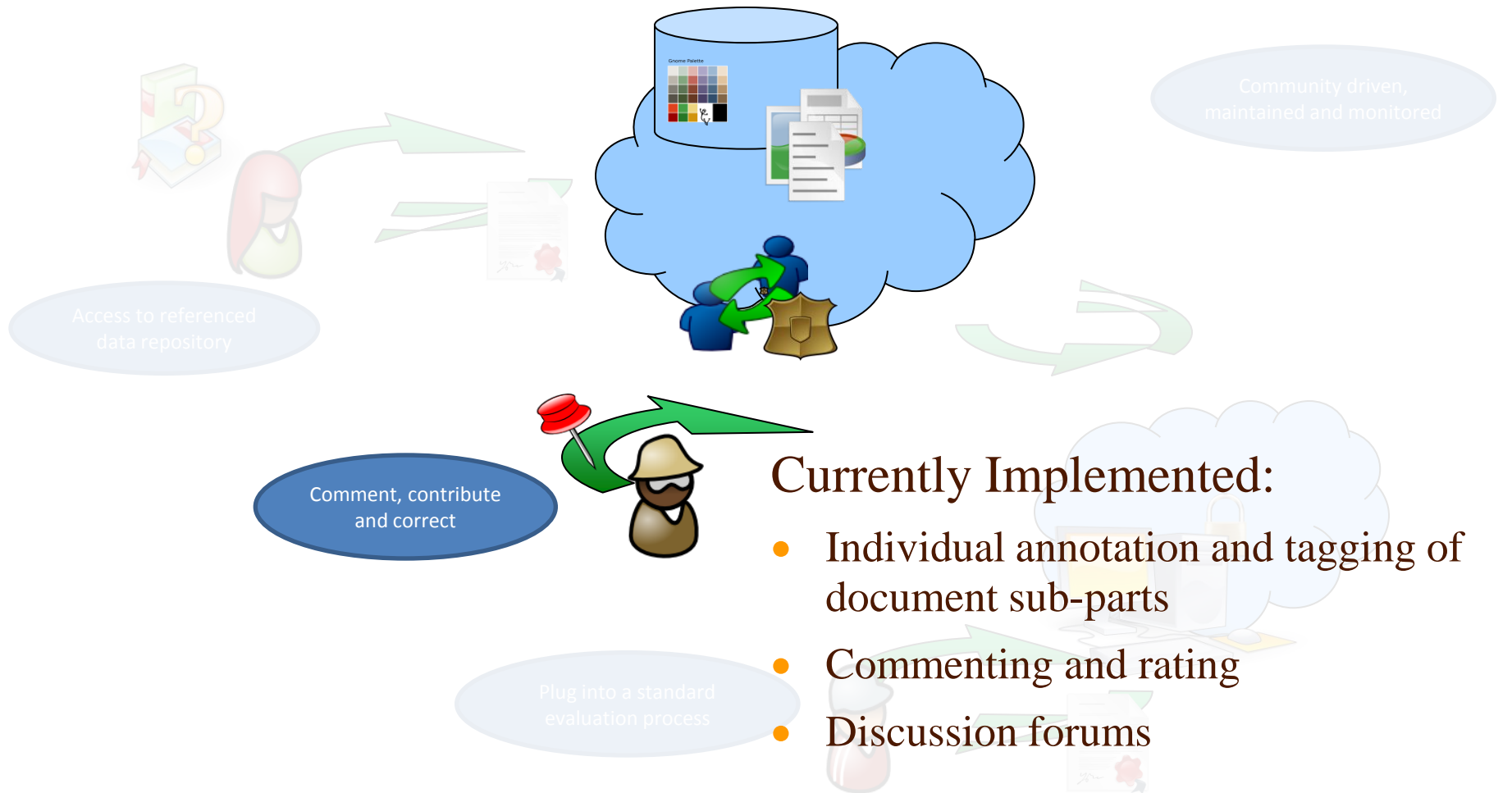
Detailed Document Description

Browse & Navigate

XML Upload



Web 2.0 Interaction





Document Analysis and Exploitation

Currently hosting
43,508 Document Images (8.9G)
15 Algorithms
167,529 Data Items

Search this site:

admin

- Browse Data
- ▷ My Uploads
- Algorithms
- My Runs
- My account
- ▷ Create
- ▷ Admin
- Log out
- Repository

Home

1_3.JPG

Submitted By: John Applesseed



Document Tagging

Document Part Annotation

Document Analysis and Exploitation newsletter

Stay informed on our latest news!

User:
admin

[Previous issues](#)

★★★★☆ admin

This document is quite nicely annotated and shows embedded and overlapping annotations. However, bounding boxes are not very precise.

Rating & Commenting

Newsletters and Discussion Groups

CKEditor: the ID for excluding or including this element is `browse/dataitem/149394.edit-dae-data-comment`.

Running Algorithms on DAE

Algorithm Hosting

Document Analysis and Exploitation

Search this site:

admin

- Browse Data
- ▷ My Uploads
- Algorithms
- My Runs
- My account
- ▷ Create content
- ▷ Administrator
- Log out
- Repository

Access to refer data repository

Home

My Runs

convert	In Progress
test both input and output	Run Failed
Tesseract (1_9.tif)	Complete
Input Files:	Output Files:
○ 1_9.tif	○ 1_9.tif.txt
Run Statistics:	
Start Time: SEP 13, 2010 at 02:26:34 PM	
End Time: SEP 13, 2010 at 02:26:53 PM	
Run Length: 0 days, 00:00:19	
Tesseract (2.tif)	Complete
Tesseract	Run Failed
ocrad (2.pbm)	Complete
ocrad	Run Failed
convert (1_7.jpg)	Complete
Tesseract (1_9.tif)	Complete
convert (1_9.png)	Complete

Home Goals People Wiki

Currently hosting
43,508 Document Images (8.9)
15 Algorithms
167,529 Data Items

Document Analysis and Exploitation newsletter

Stay informed on our latest news!

User:
admin

[Previous issues](#)



Output Retrieval

Algorithm Execution



Currently Being Implemented

Enhanced User Experience

- Improved Browsing and Retrieval
- Complex Query Interface

Resilient Queries

- Dataset Certification

Algorithm Execution Overhaul

- Full Web-based REST-WSDL Services
- Flexible Orchestration for Pipeline Execution
- Full OS Independent Virtualization

Community driven,
maintained and monitored

Access to referenced
data repository

Comment, contribute
and correct

Plug into a
validation
evaluation process



Getting Involved

Freely Downloadable Source Code (GPL)

- Standalone version available (database part still a bit tricky)
- FusionForge download site coming soon (tomorrow, promised)
- Open for contributions, help and critical user feedback



Mind Giving the
Community a Hand ?

Some Challenging Feature Requests:

- Creating a drag-and-drop query interface
- Help making algorithm execution fully web-service enabled (and thus 100% distributed and cloud-ready)
- Import / export functions to other annotation formats



Open Challenges

The screenshot shows the Google Analytics 'Create Advanced Segment' interface. The top navigation bar includes the Google Analytics logo, user information (bart.lamiroy@gmail.com), and links for Settings, My Account, Help, and Sign Out. The breadcrumb trail reads 'Manage Advanced Segments > Create Advanced Segment'. The main area is divided into two sections: 'Out of a total of ? visits...' and '...this segment matches ? visits'. The first section contains two criteria: 'Page Title' with the condition 'Matches exactly' and 'Event Action' with the condition 'Starts with'. The second section is currently empty. A 'Name segment:' field is at the bottom. Annotations include a blue box 'Construct Queries' pointing to the criteria fields, a blue box 'Intuitive Drag & Drop Query Interface' pointing to the bottom of the criteria area, and a red box 'Mind Giving the Community a Hand?' pointing to the 'Content' list on the left. A 'Test Segment' button is present in both sections.

Google Analytics

Analytics Settings | View Reports: dae.cse.lehigh.edu | My Analytics Accounts: DAE Website

Manage Advanced Segments > Create Advanced Segment

type to filter

list view

Select Criteria

Content

- Page Title
- Host Name
- Page
- Site Search Status
- Event Category
- Event Action
- Event Label

Out of a total of ? visits...

Construct Queries

Condition Value

Page Title Matches exactly

or

Event Action Starts with

case sensitive

or

Add "or" statement

and

Add "and" statement

...this segment matches ? visits

Intuitive Drag & Drop Query Interface

Name segment:

Visible in: dae.cse.lehigh.edu and 0 other profiles.

Test Segment

Test Segment

Drag and drop dimensions and metrics into the boxes to create a visit segment.

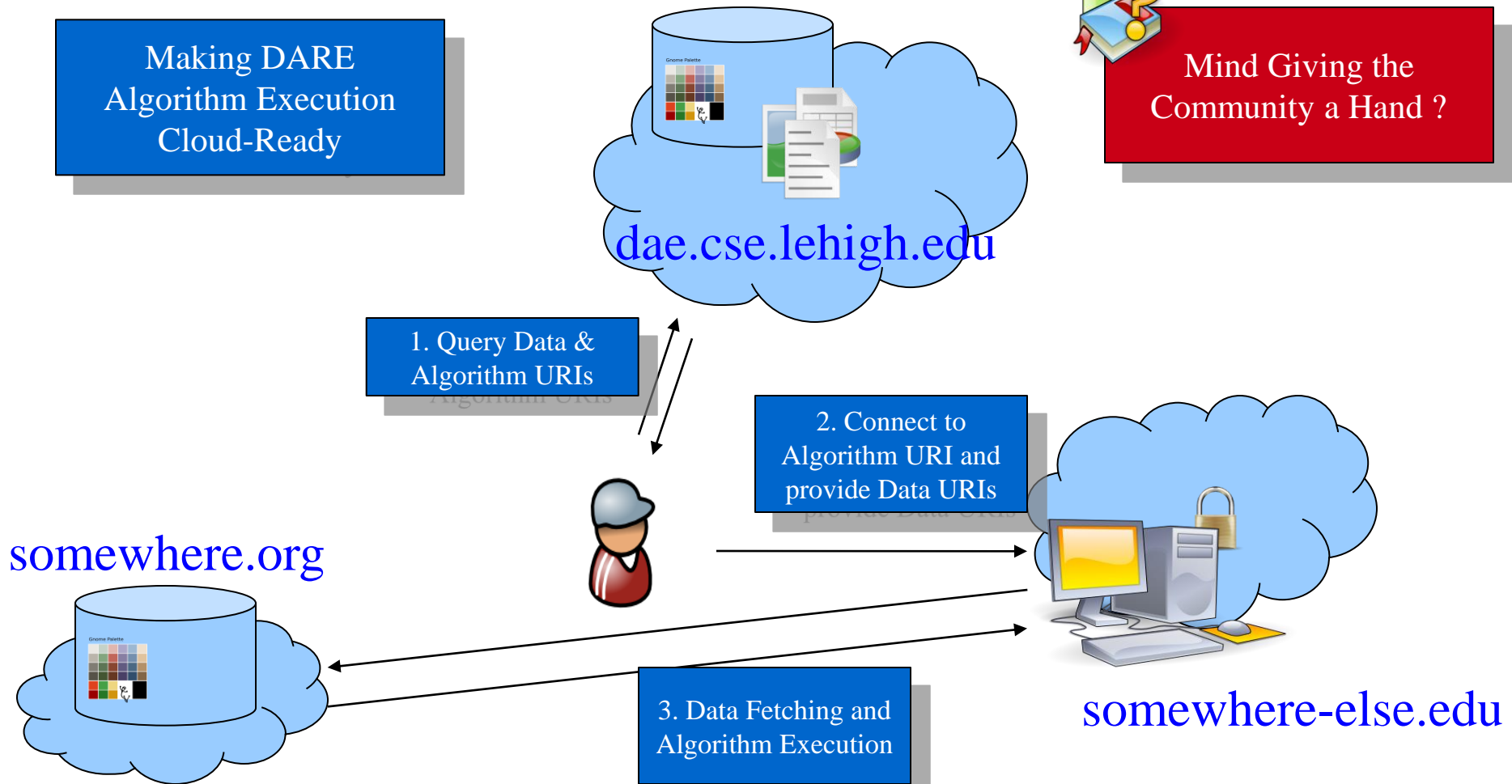


Mind Giving the Community a Hand ?

Open Challenges

Making DARE
Algorithm Execution
Cloud-Ready

Mind Giving the
Community a Hand ?



Further Challenges

- Guaranteeing Reproducible Queries
- Controlling Data Clutter
- Implementing a Distributed Data Model
- Maintaining Data Coherence
- Batch Data Processing and Performance
- Monitoring, Tracking and Exploiting Reputation



Add Your Thoughts to DARE Wiki

Contents [hide]

- 1 The Dare Wiki - What is All This About ?
- 2 Initial Thoughts and Ideas
- 3 How to Contribute to this Wiki
 - 3.1 Concept Part
 - 3.2 Scenario Example Part
 - 3.3 Implementations Part
 - 3.4 Scientific Issues Part
 - 3.5 References
- 4 Documents and Resources
- 5 Financing and Promoting Dare
- 6 References

The Dare Wiki - What is All This About ?

The Dare Wiki is a crazy-creative-innovative-visionary-insane (check the one that suits you best) idea to develop and promote the [Dare Paradigm](#).

Since the [Dare Paradigm](#) itself is fundamentally community driven, it seemed obvious that its own definition should be sprouting from and evolving through community efforts and contributions. And this is it ! Through this Wiki, we have thrown our initial ideas and thoughts to you, the community, believing you will catch up on them, enrich them, challenge them, destroy them, rebuild them, reshape them so that, eventually, given all these contributions and collective thoughts, they start living on their own.

From now on, we don't own this anymore. Connect and modify, comment or amend anything you consider appropriate.

A New Paradigm for Research in Machine Perception

(including Document Analysis and Exploitation*)

Daniel Lopresti
Computer Science & Engineering
Lehigh University
Bethlehem, PA, USA

Bart Lamiroy
Nancy Université
Campus Scientifique, BP 239, 54506
Vandoeuvre Cedex, France

* Supported by a Congressional appropriation administered through DARPA IPTO via Raytheon BBN Technologies.

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



Acknowledgements

Technical Contributors:

- Dezhao Song
- Mike Kot
- Weidong Chen
- Yingjie Li
- Yang Yu
- Xingjian Zhang
- Mike Caffrey
- Sai Lu Mon Aung

Lehigh Faculty:

- Prof. Hank Korth
- Prof. Jeff Heflin
- Prof. Brian Davison

Supporters to date:

- Prof. G. Fink (Technische Universität, Dortmund, Germany)
- Dr. D. Karatzas (Universitat Autònoma de Barcelona, Spain)
- Prof. K. Kise (Osaka Prefecture University, Japan)
- Prof. J. Lladós (Universitat Autònoma de Barcelona, Spain)
- Prof. G. Nagy (Rensselaer Polytechnic Institute, USA)
- Prof. J.M. Ogier (Université de La Rochelle, France)
- Prof. K. Tombre (INRIA, France)

