# PUBLISH ON DEMAND

# mediaX
## STANFORD UNIVERSITY

**mediaX at Stanford University**

mediaX at Stanford University connects businesses with Stanford University's world-renowned faculty to study new ways for people and technology to intersect.

We are the industry affiliate program to Stanford's H-STAR Institute. We help our members explore how the thoughtful use of technology can impact a range of fields, from entertainment to learning to commerce. Together, we're researching innovative ways for people to collaborate, communicate and interact with the information, products, and industries of tomorrow.

# Human Sciences and Technologies for Publish on Demand

## Research Theme Update

March, 2013

## Acknowledgements

# TABLE OF CONTENTS

*Human Sciences and Technologies for Publish on Demand*

*Background and Introduction*

# Introduction: Publish on Demand Theme

*Human Sciences and Technologies for Publish on Demand:*
*Optimization, User Experience and Infrastructure*

The world of media and content is experiencing an explosion of innovation that includes how content is created, consumed and curated. In publishing, this innovation has erupted in what some call the "mass amateurization" of media and extends to how traditional content creators and distributors are restructuring around new business models. One can observe escalating desires for personalization, customization, portability, and screen-of-preference viewing.  One can observe the steady stream of new devices, platforms, mobile apps and services that blanket the media landscape. One can observe new data and new data flows in the growth of self-publishing and open publishing (in educational, research, trade, and leisure content). Across media ecosystems in both education and entertainment, new technologies and new uses of them are creating a "sea change" in publishing. Interdisciplinary approaches are needed to understand how overlapping factors will influence publish-on-demand. And to explore "signals of change," structural dynamics, legal implications, and user experiences in publish-on-demand for learning, work and leisure.

Open publication/education, 'unbundling' content, and legal automation for easy rights registration and clearance are potentially disruptive technologies influencing context and format (for attention, retention, and impact of content).  New technical, human and legal issues influence analytics and pricing models in publish-on-demand systems. This influence includes choice-based options and advertising, or sponsorships. It also includes the fast-changing relationships and dependencies in the business ecosystems surrounding content in higher education. New technology developments in querying languages, linking schema, convertible formats (print, PDF, HTML), and content delivery systems with smart copyright solutions provide promising alternatives for transformation from the legacy printing industry to digital, publish-on-demand.

Six Stanford projects are exploring new insights that will help optimize user experiences and business infrastructures in the publish-on-demand mediascape of the future. This mediaX Research Theme Update provides a mid-year window into these projects. Each project is unique to this theme, yet nests into other portfolios of research activities underway at Stanford and beyond. At its core, each project builds on prior research and is intended to serve as a foundation for future research.  The projects included in the Publish on Demand Research Theme span elementary and secondary education to higher education and also include scholarly work and consumer content.  These projects were launched in September 2012 (for the 2013 academic year).

**The Transparent Social Footprints project** (co-led by Jeffrey Heer, Ann Grimes, and Jay Borenstein), is examining requirements and developing prototypes for metrics and tools to support media organizations in the shift from print to digital content. Two teams in an upper level two-quarter computer science course (CS201) are working on a use case involving a mid-sized, regional newspaper.

**The Content on the Go project** (led by Ramesh Johari) is taking a combined data-driven and structural analysis approach to study the relationship between pricing decisions and marketplace visibility. The team focuses on the role of rankings and recommendations in driving demand; they distinguish between indirect effects, such as "top rank" established through the "wisdom of the crowds," and direct effects, such as sales rank. Early results using their stylized dynamic model validate the impact of indirect effects and indicate price cycles that induce variation in rank position.

**The TweakCorps project** (led by Scott Klemmer) is investigating a hybrid machine-learning system to support the adaptation of webpages to users' needs and preferences. Building on earlier research that used an automated process to classify webpage components, the team has used expert ratings as input training data for machine-learning. The resulting automated system then leverages human decision-making to improve classification and adaptation; and to optimize the effectiveness of a user interface for editing webpages that are automatically adapted for a wide variety of user and device requirements.

**The Decision Products and Long-Term Integrity project** (co-led by Robert Laughlin and Neil Jacobstein) is assessing the human and technical infrastructures required to support the systematic construction of complex group decision products, and the procedures necessary to ensure their long-term integrity. The team is studying requirements for integrating vast amounts of data, making sense of it, and producing decision products that persist. They are also exploring the economic incentives necessary to preserve the quality, accountability and longevity of such products.

**The Smarter Scholarly Texts project** (led by John Willinsky) is exploring requirements for the use of XML (extensible markup language) for digital transformation of text-intense scholarly work. The Stanford team, in conjunction with collaborators at Simon Fraser University, the University of Heidelberg, the University of Manchester, and the University of Chicago, is prototyping a functional application that will reduce the time spent on manual editing: A multi-purpose indexing with an assisted automated system for document markup for multiple devices (including print-on-demand). The functionality includes reference checking, a parsing engine, and copy-editing. All of this will be integrated with the Open Journal Systems for multiple audiences.

**The Recasting the Textbook project** (co-led by Sam Wineberg and Roy Pea) is exploring the transformation of the textbook as an on-demand, collaborative collection of historical narratives. This digital textbook draws on primary source materials culled by high school students (from national archives, local libraries, and potentially the photo albums and historical records of students and their families). The team is also investigating the potential use of near-field communication tags and QR codes for smart phones and mobile devices to allow students to participate and interact with out-of-school resources (such as libraries, museums and social communities).

These projects continue the mediaX Publish on Demand Research Theme and its broad insight agenda on the question: *What insights about people and technology are needed to ride the sea change of publish-on-demand into the future?*

Projects sponsored earlier by mediaX under the Publish on Demand Research Theme provided a foundation of legal and infrastructure insights, which inspired exploration into prototype requirements for the Stanford Intellectual Property Exchange (SIPX). Those investigations involved expertise from computer science, law, education, business, and library information sciences. They involved the collaboration of professors, students and administrators at Stanford University, as well as Stanford University bookstore. They required the participation of service providers, decision makers and engineers in businesses whose products and services are relevant to the new publish-on-demand ecosystem. The insights from those investigations resolved some questions, reframed some questions, and raised others. This is an expected outcome for frontier research. The current Publish on Demand Research Theme continues the exploration.

**Martha G Russell**

martha.russell@stanford.edu

Executive Director

## mediaX Team

**Susana Montes**

susanam@stanford.edu

Communications Manager

**Adelaide Dawes**

adelaide@stanford.edu

Program Manager

# Transparent Social Footprints

*A New Road to Digital Dollars?*

---

*Research Team:* Jeff Heer, Assistant Professor of Computer Science; Ann Grimes, Lorry I. Lokey Professor of the Practice/Department of Communication, and Director Graduate Program in Journalism; Jay Borenstein, Lecturer Computer Science; R.B. Brenner, Lecturer, Department of Communication Graduate Program in Journalism.

*Student Researchers:* Team A - Anna Li, Danielle Radin, Xiaohuo Cui (Journalism), Michael Christensen-Calvin, Nicholas James Latourette (Computer Science); Team B - Rachel Estabrook, Ian Jacob, Riva Gold (Journalism), Azmaan Onies, Gavin Bird, Hugh Cunningham (Computer Science).

---

## Project Objectives

This project is examining how media organizations might shift from the digital metrics used by most newspapers – heavily reliant on unique visitors and pages views – to systems that better track user engagement and the behavior of individual users on the web and mobile devices. This model of a personal digital footprint, which is old hat to companies like Amazon and Facebook, could have profound implications for the delivery of content and advertising.

## Course Format

Two interdisciplinary research teams of eleven journalism and computers science students are pursuing this project as part of CS 210, an upper-level computer science course sequence in which Stanford student teams collaborate on software challenges that require innovation. Teams often work on projects funded by corporate sponsors. In this case, they are tackling the Transparent Social Footprints project. In this course, student researchers take projects all the way from concept to completion. This includes defining requirements, iterating through ideas and prototypes and, ultimately, producing a prototype for a software product. Projects span all industries and focus on applying next generation software approaches and techniques to yield better solutions.

## Use Case

The use case selected for this project is the McClatchy Corp., which operates 30 newspapers in 15 states. The team is targeting The Sacramento Bee, which has been identified as "typical" for a mid-tier newspaper chain. There are 1,400 daily newspapers in the U.S. and the team believes the data management problem of tracking and monetizing user engagement is common for many U.S. news organizations. The team also selected the Bee because McClatchy agreed to share the paper's back-end databases for research purposes. The Bee is located within a reasonable geographic distance so students and company liaisons can meet both on campus and on site in Sacramento.

## Design Parameters and Targets

The team is tackling this problem from two angles; both are open source, for broad industry use. Since the start of winter quarter, the student teams have moved from concept to analysis to initial design to user testing. They are now coding and will have initial software prototypes ready by March 18. The student researchers also have met with fellow mediaX grantee Prof. John Willinsky and his team to brainstorm on potential synergies with his group.

## Team A - Project Results to Date

The team is finishing the initial version of an algorithm that builds up a set of metadata when it doesn't exist, and an API for publishers that supports existing metadata. The algorithm identifies certain keywords that define a piece of content, making it easier to match an individual reader's interests with other content – text, videos, advertisements, etc. – likely to appeal to him or her. The team is using MongoDB, a document database, and tapping into a mix of semantic and social media tools – including Open Calais, Lingospot, Gigya and Chartbeat – to harvest metadata and better link it to the user. Users are being tracked through a consolidation of online registration tools: Facebook and other social networks; free and paid web registrations; print subscriptions; and a range of apps that ask for various personalized identifiers, including geo-coding.

## Team B - Project Results to Date

This team is focusing on heightened reader engagement. The team is motivated, in part, by Facebook's proposition that content with robust user activity around it is of greater value-add than, say, a story link in the news feed that yields no comments or few comments. The team has completed an initial prototype that enables website and mobile users, when reading an article, to access an overlay of relevant information, as well as options for participating in discussions. The team envisions several use cases. These allow for the identification of "hot topics" that track the digital footprints of "fans" (the most loyal readers) based on the kinds of stories they are connecting to, thereby distinguishing between more and less valuable digital readers. User testing is now underway to further refine this prototype.

## Spring Quater Targets

Student prototypes will serve as the basis for further user testing and product refinement. Student teams will present their prototypes to publishers starting spring break week of March 25 for further feedback. In addition to McClatchy, the groups will present to: The New York Times, The Wall Street Journal Digital, Chartbeat, Open Calais (ThomsonReuters), DigitalFirstMedia and, during spring quarter, Konica Minolta's representatives. Based on feedback, the week of April 1, teams will reset spring quarter weekly milestones. The team anticipates projects to be completed by June 7, with next-stage roadmaps drawn up at that point.

# Content on the Go

*The Economics of the Market for Mobile Apps*

*Research Team:*  Ramesh Johari, Associate Professor Department of Management, Science & Engineering; Bar Ifrach, Postdoctoral Research Fellow, Department of Management Science & Engineering.

## Motivation and Introduction

The growth of mobile applications on smartphones and tablets ("apps") ranks as one of the most astonishing technological developments in recent past. Over 600,000 apps, either free or paid, are available for immediate download from designated app markets (e.g., App Store and Google Play). These app marketplaces are a significant disruptive change in the way content is created and consumed. On the demand side, the marketplaces offer users rich content utilizing the functionality of their mobile devices, grabbing their attention away from legacy media (e.g., print) and even modern media (e.g., web browsing on a desktop). On the supply side, these platforms provide content creators direct, instantaneous, and highly popular distribution systems where they can implement their own marketing and pricing policies, cutting out middlemen. However, for a content creator, making sensible business decisions requires an understanding of the economics underlying this market---including competition, features of the market platform, and pricing.

Taking a combined data-driven and structural analysis approach, this project studies various aspects of the relationship between pricing decisions and marketplace visibility.  The team's aim is to empower individual content creators by offering strategic guidance on how to leverage the marketplaces' flexibility. In particular, the team focuses on the role of rankings and recommendations in driving demand. The market platforms offer a number of recommendation systems designed to harness the so-called "wisdom of the crowds" to help users choose what to download in the plethora of apps. The most salient among them are the "top-ranked" charts that list apps by number of downloads, as well as some secondary popularity indicators. A high position in these charts is followed by a remarkable boost in demand, according to both industry sources (Surikate and GfK, 2012) and empirical research (Carare, 2012). The team calls the effect of top-rank position on future sales an indirect effect, to distinguish it from the direct relationship between the past sales and rank (since rank is a measure of past sales expressed in comparison to those of competing apps).

The team proceeds in two directions. First, the team postulates a reduced form model to estimate the magnitude of the indirect effect of ranks on sales, employing time series of top-ranked charts in the second half of 2012. This model isolates the indirect effect, as outlined in more detail below. The team's results show that the indirect effect is statistically significant and substantial.

With this effect in hand, the team studies app pricing decisions in a stylized model that incorporates earlier findings. Surprisingly, the team finds that accounting for the indirect effect may give rise to optimal price cycles, in which the seller alternates periodically between a high price and a low one to boost revenue in the first, and market visibility in the latter. The team finds evidence in the data supporting this pricing behavior in practice.

In the reminder of this abstract, the team provides more detail on the empirical and structural components of this project, and reports some of the findings.

## Empirical Model

Top-ranked lists capture the popularity of different apps based on their recent demand, rewarding popular apps with salient market visibility and a trendy appeal. However, as much as ranking lists reflect underlying demand patterns, they may also set them by increasing the demand of already top-ranked apps over less popular ones, further skewing the demand distribution of apps. In the presence of this indirect feedback, the demand for a particular app will be a function of not only its attributes (e.g., functionality, graphics, description, etc.) and price, but also of its rank position.

Ranks change over time for two reasons: (a) reflecting organic variations in the demand for the app itself, and (b) reflecting exogenous variations in the demand for other apps. For example, between July 11-12, 2012, a number of highly ranked apps dropped one position in the overall top-paid list, following the launch of the much anticipated game *Amazing Alex*, which went straight to the top of the chart (the game later flopped). A popular texting app, *WhatsApp Messenger*, dropped from position 3 to 4, but this decrease in rank was likely unrelated to a change in its underlying demand.

On the other hand, on the same dates, a photography app, *Camera+*, dropped from position 7 to 9; while the rank of an emoticon app, *Emoji 2*, jumped up from rank 9 to 8. The team calls this a "swap": the ranking order of two apps on the same list changed. A swap indicates an organic change in demand for at least one of the swapping apps: either the demand of the initially less popular app increases or that of the more popular app decreases (or both). The team uses this distinction to isolate the indirect effect of ranks on demand. In particular, the team works from the principle that a variation in download rank that does not include a swap reflects an exogenous change in rank. This approach is used as a basis to estimate the indirect effect of rank on future sales.

The team focuses on the most visible ranking lists: the overall (i.e., across all categories) top-free, top-paid and top-grossing. The first two capture download popularity for free apps and for apps with a positive download price, respectively. The grossing ranking captures the revenue generated by apps, both from the download price and from in-app purchases (e.g., upgrades and additional functionality), and includes both free and non-free apps. Although the platforms do not disclose the algorithm behind these ranks, the common belief in the industry is that the algorithm is primarily based on a weighted average of recent downloads (top-free/paid) or revenue (top-grossing). The team's dataset includes a time series of these charts, recorded 2-4 times daily. The team makes use of the fact that many apps appear in both top-paid and top-grossing lists, and their econometric model relates variations in top-paid ranks to variations in top-grossing ranks that are not the result of a swap.

The team's results establish that the indirect effect exists and is significant for almost all relevant rank ranges. To illustrate, the team finds that the indirect effect of dropping one position among the top 5 positions in the top-paid chart is a 5.4% decrease in revenues on average (95% confidence interval is 3.4-7%).[1]

Dropping from the top position to the fifth in the top-paid chart would result in an impactful 21.6%

---

1 This estimate uses the result in Garg and Telang 2012 to relate the rank position to revenue.

decrease in revenues. Additional regularities in the team's estimates further support their findings.

## Structural Model and Price Cycles

Motivated by the significant indirect effect of rank position on sales, the team proposes a stylized dynamic model for a developer seeking to maximize her expected discounted revenues in the marketplace. To account for the indirect effect, the team supposes that sales at the present period are positively affected by those in the previous one. The developer can adjust her price dynamically, under the standard assumption of an inverse relationship between price and sales (downloads). Surprisingly, the team finds that this simple deterministic setting can give rise to optimal price cycles. Intuitively, the developer alternates between boosting next period's demand by dropping the price in the current period, and monetizing on that with a higher price in the following period. These results are related to a stream of literature in economics studying optimal cycles and chaos in growth models to explain business cycles (See Nishimura and Sorger, 1996).

Examining developers' pricing policies in the team's dataset, they find a large number of developers using price cycles as suggested by the structural model. In addition, the team's discussions with industry practitioners indicate that in many cases these cycles are aimed to induce variation in the rank position to improve marketplace visibility. The team is currently collaborating with developers to further study the ranking mechanism and test these predictions.

## References

Carare, O. (2012) the Impact of Bestseller Rank on Demand: Evidence from the App Market. International Economic Review 53, 717-742.

Garg, R. and Telang, R. (2012) Estimating App Demand from Publicly Available Data. Mimeo.

Nishimura, K. and Sorger, G. (1996) Optimal Cycles and Chaos: A Survey. Studies in Nonlinear Dynamics & Econometrics 1, 3.

Surikate and GfK (2012) An insight into iPhone user behavior within the App Store. URL http//www.surikate.com/en/etudes.html

# TweakCorps

## *Re-targeting Existing Webpages for Diverse Devices and Users*

Research Team:  Scott Klemmer, Associate Professor, Computer Science; Maxine Lim, Graduate Student, Computer Science.
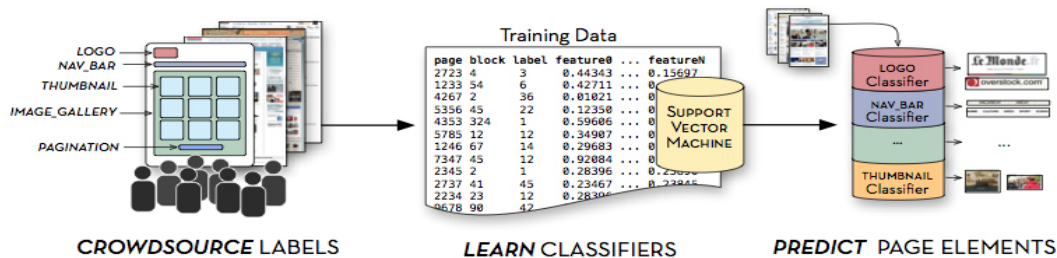


*Figure 1:  The pipeline for learning structural semantic classifiers for the Web. First, a large set of labeled page elements are collected from online workers. Next, these labels are used to train a set of regularized support vector classification SVMs. These classifiers are then used to identify semantic elements in new pages.*

## Project Overview

The research team has developed a machine-learning algorithm that can label web pages. The algorithm takes the structure of pages, conducts design mining and provides the opportunity to learn about design elements from multiple sources.

While designers would prefer to design for as few devices as possible, the incentive of hardware manufacturers is to offer a "new market." The team worked on this learning algorithm to solve the tension between hardware manufactures, who seek to develop new systems, and designers, who prefer to not make too many changes to their web design work when new devices are developed.

The team's algorithm enables a new kind of design-based machine learning, that can stream structured visual descriptors for page elements from a central repository. The algorithm also allows data to be collected and integrated with the repository for supervised learning applications like crowd-sourcing.

The algorithm takes the structure of the web pages, and analyzes them. It is composed of five integrated components:
- A web crawler
- A proxy server
- The data store
- The post-process
- The API

For this project, the team focused on one popular class of semantic identifiers: Those concerned with structure – or information architecture – of a page. The team explored a different tactic for adding structural semantics to web pages; with accurate learning classifiers, pages can be semantified automatically, in a post-hoc fashion, decoupled from the design and authoring process. To this end, the team presents a classification method based on support vector machines (i.e., a supervised learning model with algorithms that analyze data and recognize patterns), trained on a large collection of human-labeled page elements and employing a feature space comprised of visual, structural, and rendered-time page properties (Lim, Kumar, Torres, Talton, Satyanarayan, Klemmer, 2013) (see Figure 1).

The team took a crowd-sourced approach and recruited 400 participants on Amazon's Mechanical Turk, who collected more than 21,000 semantic labels over a corpus of over 1,400 web pages. The team used labels to determine the set of classifier and provide training data for machine learning (Lim, Kumar, Torres, Talton, Satyanarayan, Klemmer, 2013). Every participant applied semantic labels to at least 10 elements on each of five pages. The pages used in the study were drawn from the Webzeitgeist design repository, which provides visual segmentations and page features for more than 100,000 webpages. Webzeitgeist was a platform developed by the team to help users understand design demographics, automate design curation, and support new data-driven design interactions (See: http://vis.stanford.edu/papers/webzeitgeist)

To more thoroughly understand how labels relate to one another, the team created a co-currence matrix for the 85 most-frequent labels, each of which was used 20 or more times. To evaluate the feasibility of learning structural semantics from data, the team trained binary SVM classifiers for the study's 40 most frequent labels.

## Challenges

•       Webpages have different elements and are made by different people. There is no consistency in data. This raises the question: How can you deal with unruly data?

•       On the web, the main challenge is building systems that can deal with messy, badly formatted data.

•       Designs have more and different types of structures.

## Next Steps

•       Maxine Lin will submit the "Learning Structural Semantics for the Web" paper and will present a live demo of the prototype's search capability.

•       The team would also like to get a preliminary version of a responsive design (short-term horizon).

•       Continuing work that was been underway for 2 1/2 years, the team plans to improve accuracy of the classifiers and have search on Webzeitgeist ready by the summer.

# Publish on Demand Decision Products

## *And Their Long-term Integrity*

*Research Team:* Robert Laughlin, Anne T. and Robert M. Bass Professor of Physics; Neil Jacobstein, Research Associate and mediaX Distinguished Visiting Scholar at Stanford University.

## Background

This project is designed to assess the human and technical infrastructure required to support the systematic construction of complex group decision products, and the procedures necessary to ensure their long-term integrity. It is actively addressing two interdisciplinary research questions:

1) What infrastructure requirements would support the systematic construction of decision products that are designed with explicit requirements, grounded by open and well-documented research, subject to systematic review and quality control, and communicated persuasively to relevant constituents?

2) Given the explosion of ephemeral digital media, what are the criteria for educational, technical, and institutional mechanisms to ensure the long term integrity of decision products and intellectual work?

The rapidly increasing velocity and complexity of decision-making in industry, government, and non-profit organizations is making business-as-usual methods for situation assessment and decision-making untenable. Information of uncertain quality and accountability is exploding on the web. Our methods for deliberation in government agencies, board rooms, and non-profit groups need to change to meet the accelerating information challenges.

Specifically, the human brain evolved under very different linear and local event selection pressures than the exponential and global events in today's business environment. Further, the brain has a well documented set of hard-wired or built-in heuristics, such as overgeneralization, saliency, and sunk cost biases, that were once adaptive to keep us safe from predators, but are maladaptive distortions of judgment in the complex modern world. Even well-educated decision makers are subject to these and other systematic errors in judgment. However, it is possible to provide explicit anti-biasing decision support that decreases the likelihood of systematic decision errors.

The research team suggests that what is needed is software infrastructure for developing a culture that makes high quality decision products a priority - and preserves a record of the decision support information and its provenance.

## Progress to Date

The team is researching conditions for augmenting our limited processing capabilities for integrating vast amounts of data, making sense of it, producing decision products that persist, and feeding back the results of our current decisions into our future decisions. The team is also investigating the economic incentives necessary to preserve the quality, accountability, and longevity of these decisions.

The team designed and conducted a mediaX Workshop on Augmented Decision Systems Concepts, Incentives, and Requirements. It was held on March 2, 2013.  Distinguished technical, corporate, and government leaders were invited to help the team think through system use and requirements in the new and innovative Peter Wallenberg Learning Theater high-resolution video display wall environment at Stanford University.

The workshop generated many novel and exciting ideas. Several of the participants are interested in pursuing these ideas.

## Preliminary Workshop Results

**1.** Workshop participants agreed that technology alone would not solve the complex decision quality problems in corporations, government, and non-profit organizations. However, a principled combination of innovative decision methodology, human incentives engineering, and decision support software and hardware technology could make a significant difference in decision product quality.

**2.** One of the needs identified by workshop participants was to provide a decision team with support for communication, visualization, and touch-based manipulation of complex decision products.  For example, one of the breakout groups explored a use case that addressed the needs of a Rapid Issue Response Team that could be charged by a corporate executive or government agency to rapidly gather the relevant information about an urgent or critical issue, test facts and identify misinformation, assemble a balanced and coherent representation of "n" sides of the issue, develop recommendations, and work products to explain the recommendations to relevant audiences in clear and compelling ways.

**3.** Providing resilient infrastructure for decision-making under crisis conditions and uncertainty was another theme. The use case for this need was providing effective emergency response to a massive earthquake, in spite of the absence of grid electricity and cell phone communications.  There were many very useful ideas knit together into one powerful and surprisingly resilient emergency response decision system.

**4.** Addressing dysfunctional government decision making was another use case, that focused on the problems associated with government and corporate decision making systems that were once functional, but have become slow, self-serving, simplistic, and non competitive with best practices. There were many excellent ideas generated under this theme. One in particular that stuck, and has already attracted a team that wants to pursue it, is the idea of a "People's Accountability Engine" that could analyze complex reports and policies, and reduce the text to a set of key assertions, assumptions, and additional

action items not directly related to the purpose of the legislation or report (otherwise known as "pork" or favors for special interest groups). Information gleaned from the report could be utilized to simulate probable outcomes under varying conditions and assumptions. Stanford Professor Jim Fishkin, participated in this session. He has conducted more than 70 large-scale deliberative decision making workshops in over 15 countries. He thinks that an augmentation system like the "People's Accountability Engine" discussed in the workshop, would be a transformative technology for government and business decision-making.

**5.** The workshop highlighted the utility of sophisticated visualization displays and software tools. Rather than have a special visualization cave, which is the only place to experience substantial intellectual augmentation, the participants expressed an interest in having a seamless augmentation environment – from smart phones, to wearable displays such as Google's monocle "Glass" product, to slates and laptops, to wall sized and 320 degree immersive 3D displays.

**6.** All three groups pursuing different application use cases identified the usefulness of having a system that could do semi-automatic fact-checking, by searching the web for counterfactuals to assertions made. In addition, a system that in its simplest form could be a checklist for avoiding known bugs and biases in cognition, such as those identified by Tversky and Khaneman, could be enormously valuable if implemented rigorously. AI techniques such as the Deep Learning algorithm give us the ability to recognize these pathological patterns with much greater fidelity than simple check-lists.

**7.** Once a high quality decision product has been produced, there is still the matter of ensuring its integrity and longevity. One of the techniques for ensuring integrity is to score the information sources used to produce the product. Another method is to enforce the constraint that all references either be backed up in multiple places, or available on paper.

## Phase I - Project Objectives

**1.** Develop a decision products methodology and initial virtual decision support environment: The team is converging on a minimal methodology. Following the workshop, the team will select the most frequently cited required augmentation applications, and build them into a virtual reality for demonstration and testing.

**2.** Prototype and extend a cloud-based course that can be used to evaluate educational, technical, and institutional mechanisms required to ensure the long term integrity of intellectual work products, including decision products: Professor Laughlin prototyped an argument accountability application for use in his Stanford course. The high quality results have demonstrated the effectiveness of the process (See http://large.stanford.edu/courses/2011/ph241/). The team is targeting several potential project-specific improvements.

**3.** Run the mediaX Workshop on Decision Products and Information Integrity: This has been accomplished.  The workshop was held in the Wallenberg Learning Theater at Stanford University.

**4.** The research team plans to test open interfaces to a hybrid electronic storage and print-on-demand system for documenting decision products and student-developed course work. The team needs the API (applications programming interface) of the print-on-demand system in order to create links with the Phase I demonstration software.

**5.** Organize a mediaX Workshop on Decision Products and Information Integrity: This is done. This workshop helped the team understand what works currently, what doesn't, and what else is needed. The team is gleaning additional insights from the products of the workshop, and is drafting the related report.

**6.** Project Report: The team is working on a Phase I report that will document preliminary findings, lessons learned, a walk through of current and proposed software environments, and next steps for Phase II. The report will specifically address some exciting areas of new opportunity to provide seamless and operationally functional decision augmentation that could be a game changer for early-adopter organizations.

It now appears that in the ten year forecast, most organizational decision-making will be significantly augmented to enhance decision product quality, including considerable information integrity checking and persistence assurance. Moving proactively to develop, and especially, adopt products and services that could accomplish this will confer significant competitive advantage.

# Appendix

## Workshop Participants

**Daisuke Asai**

Researcher, NTT Cyber Solutions Laboratories, Japan

Bio: http://agelab.mit.edu/daisuke-asai

**Neal Burns**

Professor College of Communication, University of Texas at Austin

Bio: http://caps.ucsf.edu/personnel/nburns/

**Harry Blount**

Founder & CEO DISCERN

Bio: http://www.discern.com/team

**Malcolm Davies**

Director of Development, Center for Understanding Change (C4UC)

Bio: http://linkd.in/Z8rjWp

**Jim Fishkin**

Chair / Professor and Director Department of Communication and Center for Deliberative Democracy at Stanford University

Bio: http://comm.stanford.edu/faculty/fishkin/

**Gene Golovchinsky**

Senior Research Scientist, FX Palo Alto Laboratory (FXPAL)

Bio: http://www.fxpal.com/?p=gene

**Martin Haeberli**

Founder and Principal at Haeberli Associates

Bio: http://www.linkedin.com/in/haeberli

**Robert Horn**

Visiting Scholar Center for the Study of Language and Information at Stanford University. Owner, MacroVU, Inc

Bio: http://en.wikipedia.org/wiki/Robert_E._Horn

**Evan Huddleson**

Entrepreneur

Bio: http://www.linkedin.com/pub/evan-huddleson/40/215/473

**Joshua Kauffman**

Advisor and Collaborator, The Quantified Self

Bio: http://www.linkedin.com/in/joshuakauffman

**Jason Leigh**

Director, Electronic Visualization Laboratory

Professor, Computer Science, University of Illinois at Chicago

Bio: http://www.evl.uic.edu/spiff/JASON_LEIGH/Bio.html

**Ben Lenail**

Director of Business Development, Alta Devices

Bio: http://www.linkedin.com/pub/ben-lenail/0/35/361

**Tim McCormick**

Blogger

Bio: http://linkd.in/12hKDH3

**Sam Perry**

President Ascendance Ventures

Bio: http://www.linkedin.com/pub/sam-perry/0/1/74b

**Tony Seba**

Entrepreneur, Author, Lecturer at Stanford University

Bio: http://www.linkedin.com/in/tonyseba

**Russell Thomas**

Researcher, Computational Social Science, George Mason University

Bio: http://www.css.gmu.edu/?q=node/104

**Hiroshi Tomita**

President at Konica Minolta Systems Lab USA

Bio: http://linkd.in/XIt0vE

## Hosts

**Neil Jacobstein**

mediaX Distinguished Visiting Scholar

Bio: http://www.imm.org/about/jacobstein/

**Robert Laughlin**

Anne T. and Robert M. Bass Professor of Physics

Bio: https://physics.stanford.edu/people/faculty/robert-laughlin

**Martha Russell**

mediaX Executive Director

Bio: http://linkd.in/15oiyvs

## mediaX Staff

**Susana Montes**

mediaX Communications Manager

Bio: http://linkd.in/gakoaz

**Adelaide Dawes**

mediaX Program Manager

Bio: http://linkd.in/Z3zO3U

# Smarter Scholarly Texts

*For Cross-platform Publishing, Text-mining, and Indexing*

*Research Team:* John Willinsky, Khosla Family Professor, Graduate School of Education; Alex Garnet, Data Curator, Simon Fraser University Library; Juan Pablo Alperin, Researcher and Systems Developer, Public Knowledge Project.

## Project Overview

The first several weeks of work (since mid-October) have largely centred around early architecting, communicating with the University of Manchester, and hiring an external software developer. Work began in earnest on the project at the end of October 2012, after successfully hiring the developer (Damion Dooley) who will do 80% of the actual system-building and web interface development for this project.

Developing pdfx in collaboration with the University of Manchester was unexpected, but came at a very fortuitous time for the team. The Stanford team was just about to begin work on a parsing engine almost from scratch (due to dissatisfaction with existing tools). The Stanford team was able to negotiate licensing terms, which allowed them to utilize the University of Manchester's parsing engine (including its source code) as a basis for development and offer its services free of charge, despite the fact that the Manchester team is actively selling it to commercial publishers. The Stanford team is getting this service, which may become an industry standard, for free– as long as the team does not redistribute the Manchester team's code. In exchange, the University of Manchester received a feedback and testing agreement, as well as access to any of the features developed on top of their code (which so far have taken the form of XML format conversions). Based on this, the team decided to pursue a service-oriented development architecture wherein the parsing engine (utilizing pdfx) would be effectively closed-source, but could be freely utilized as a standalone web service (currently functional at http://142.58.129.113/dev/). The Stanford team will then develop plugins for other PKP platforms (notably, Open Journal Systems – OJS – and Open Monograph Press, or OMP), which would call this service, and use its parsing functionality within an existing workflow.

## Progress to Date

The team has spent a significant amount of time on improving the pipeline of tools used by this parsing engine. In January, the team focused solely on the parsing of bibliographic references. Many potential problems of parsing a document body are magnified during this process, due to even tighter formatting expectations and even more already-existing open source toolkits, which needed to be made to work together. The team eventually utilized pandoc (https://github.com/jgm/pandoc) and ParsCit (https://github.com/knmnyn/parscit) as well as some of the team's own heuristics.

The Stanford team has successfully implemented rudimentary copy-editing functionality, which checks references listed at the end of the article against those used in the article body and vice versa. There is

also an elegant find-as-you-type solution for selecting a preferred citation style from the Citation Style Language repository (https://github.com/citation-style-language/styles), which ensures that the service can cater to virtually any available journal guidelines, without needing to invest any future effort into implementing new styles or style revisions.

The team successfully met a mid-January milestone with the first live demo of this service at the mediaX2013 conference. After receiving positive feedback, the team planned to focus on improving the parsing heuristics in advance of the second project milestone: the upcoming Beyond the PDF workshop in mid-March in Amsterdam (http://www.force11.org/beyondthepdf2).

However, due to delays involved in getting a workable version of the pdfx source, the team has instead turned its focus on the OJS plugin in advance of this March milestone, with parsing improvements to follow in April. Working on OJS has necessitated bringing Damion Dooley – a seasoned web developer, though unfamiliar with PKP platforms and the scholarly publishing context – into closer contact with the rest of the PKP team. This has delayed progress, but has been beneficial in terms of knowledge transfer.

On March 4th, the team had a standalone service, which produces a 75% desirable result on article body text and a 95% desirable result on article references. The team is less than a month away from having a functional, mediated interaction between Open Journal Systems' workflow and the parsing service. The team has also been in discussions with the University of Chicago, who are developing a similar system, which uses NLM XML as a back-end. However, this system generates ePub-format ebooks, which the team expects to be able to use in tandem with their own code.

The team has assisted with the preparation of a grant at the University of Heidelberg which, if successful, will provide funding for the development of a WYSIWYG XML editing system. This system can be seamlessly invoked at the end of the team's automated pipeline, both minimizing the time spent on manual editing and allowing the result of their automated system to receive the human attention, likely necessary to turn a 95% satisfactory result into production-quality output.

## Next Steps

Following the release of the OJS plugin, the team will have several remaining development priorities for the remainder of this initial funding period through May/June 2013, including:

•       Improve pdfx's parsing heuristics, particularly in the realm of multilingual support, by leveraging PKP's existing community of OJS translators. This is something which pdfx has said they will not be supporting in the short term but will be very important and achievable for the team's user community.

•       Begin testing the service in OJS production workflows, in collaboration with some of the  more prominent journals, so that the team can assess its real-world usability, and begin gathering feedback in advance of future revisions.

•       Develop a plugin for integration into Open Monograph Press.

# Recasting the Textbook

*An On-demand, Collaborative Collection of Historical Narratives Through Primary Documents and Interactive, Touch-based Devices*

*Research Team:* Sam Wineburg, Margaret Jacks Professor of Education and (by courtesy) of History; Roy Pea, David Jacks Professor of Education and Learning Sciences; Laura Moorhead, PhD candidate, Learning Science and Technology, Graduate School of Education; Molly Bullock, PhD candidate, Learning Science and Technology, Graduate School of Education; Paul Franz, PhD candidate, Learning Science and Technology, Graduate School of Education; Jeremy Jimenez, PhD candidate in International and Comparative Education; Max Alexander, MA candidate Learning, Design, and Technology.

## Progress to Date

The team's effort to recast the history textbook as an on-demand, collaborative collection of historical narratives through primary documents is well underway. During the fall, the project team expanded to include: Max Alexander, a Learning, Design, and Technology master's student with a background in instructional design and classroom technology integration; and Jeremy Jimenez, a PhD candidate in International and Comparative Education. Rob Lucas, a postdoctoral student experienced in history education and the learning sciences, joined the team as an advisor. Additionally, a master's student in human computer interaction design will join the team during the spring quarter.

During late fall, Stanford's IRB approved the research effort. The team has also secured the participants and site: A collaborating teacher and two classes of approximately 50 high-school juniors in a northern California urban charter school. On-site research at the school will begin in earnest during April and May.

## Project Objectives

The objectives around this research project remain focused on exploring how the digital textbook might be recast by students to include a variety of primary source materials culled from national archives, local libraries, and potentially the photo albums and historical records of students and their families. As part of the project, students will design and create a digital textbook about world history from Greek and Roman time through the post-World War II era. The project encourages students to critically read and judge primary sources, as well as critique and construct historical accounts that thicken the narrative through the inclusion of multiple perspectives.

An archivist from the Hoover Archives has agreed to visit the school, and students, in turn, may visit the Hoover Archives at part of the project. Additionally, a potential collaboration with archivists at the National Archives appears promising.

The team is designing the project to allow debate around each document. As an expanding on-demand digital textbook, students will be able to "rate" and comment on the historical value and validity of a source. They will be expected to create, explain, and justify multiple historical narratives using iBook Author with templates designed and coded specifically for this project.

During April and May, students and their teachers will be interviewed, videotaped, and surveyed about their digital textbook experience. The research team is developing the protocols for this effort, and the plan is to pilot these protocols over the next several weeks. Data logs, as well as each student's selected or uploaded documents and created historical narrative, will also be studied through qualitative and quantitative methods.

This research project will also explore the boundaries of on-demand, digital publishing through the use of novel technology, notably near-field communication (NFC) tags and Quick Response (QR) codes. Questions for exploration include: How might these technologies and other methods of collection and curation for smartphones and various mobile devices allow students to tap into a historical narrative at their local libraries, archives, and museums? How might these technologies encourage the discussion and "rating" of documents that allow for a historical event to have multiple "truths" or narratives, which include a variety of media? How might these technologies further push the boundaries of on-demand publishing? What can we learn about the creation process of collaborative, technology-supported media from collecting log data?

The team is also in the process of collecting and digitizing (when need be) approximately 200 historical assets (all free from copyright restrictions). Students can pull from and add to this collection as they develop their digital textbook historical narrative. Finally, a teacher curriculum to accompany the digital textbook project in the classroom is underway.

The team is working at a school that is known to be progressive, where students are asked to engage in a synthetic and creative assessment task. However, even at this site, media creation and social learning affordances of new digital technologies are not leveraged. Rarely do students dive deep into substantial, authentic sources. Much of the team's challenge going forward will be addressing this issue. For instance, how might the technological tools currently available encourage students to use primary sources to "thicken," or open up another side of historical knowing? How might students reveal the multiple and often nuanced answers to a seemingly simple question? (e.g., Why didn't Rosa Parks give up her seat on the bus?)
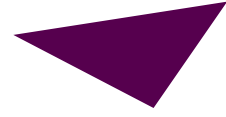
## Next Steps

Future research questions and potential projects have emerged. For example:

- How might students tasked with writing a primary-source document counterbalance the use of their textbook over the course of a school year?

- How do the technical tools currently available in school settings either encourage or limit this effort?

•       What other tools might be of value?

•       What are the roles and potential impact on students' motivation and empowerment in the production of digital textbooks?

•       What can be done in the medium to increase a student's sophistication in assessing sources and historical understanding?

# Learning Structural Semantics for the Web[1]

*Maxine Lim, Ranjitha Kumar, Arvind Satyanarayan, Cesar Torres,
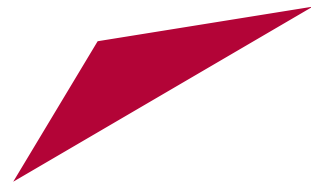Jerry O. Talton, Scott R. Klemmer.*

**March, 2013**

**mediax.stanford.edu**

The format of this paper is optimized for double-sided printing.