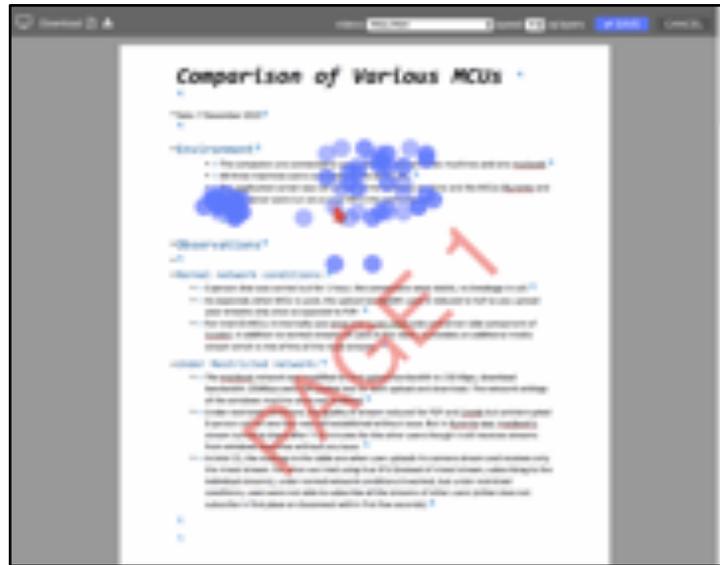


DocuGram: Turning Screen Recordings into Documents



Laurent Denoue, Scott Carter, Matthew Cooper

FXPAL

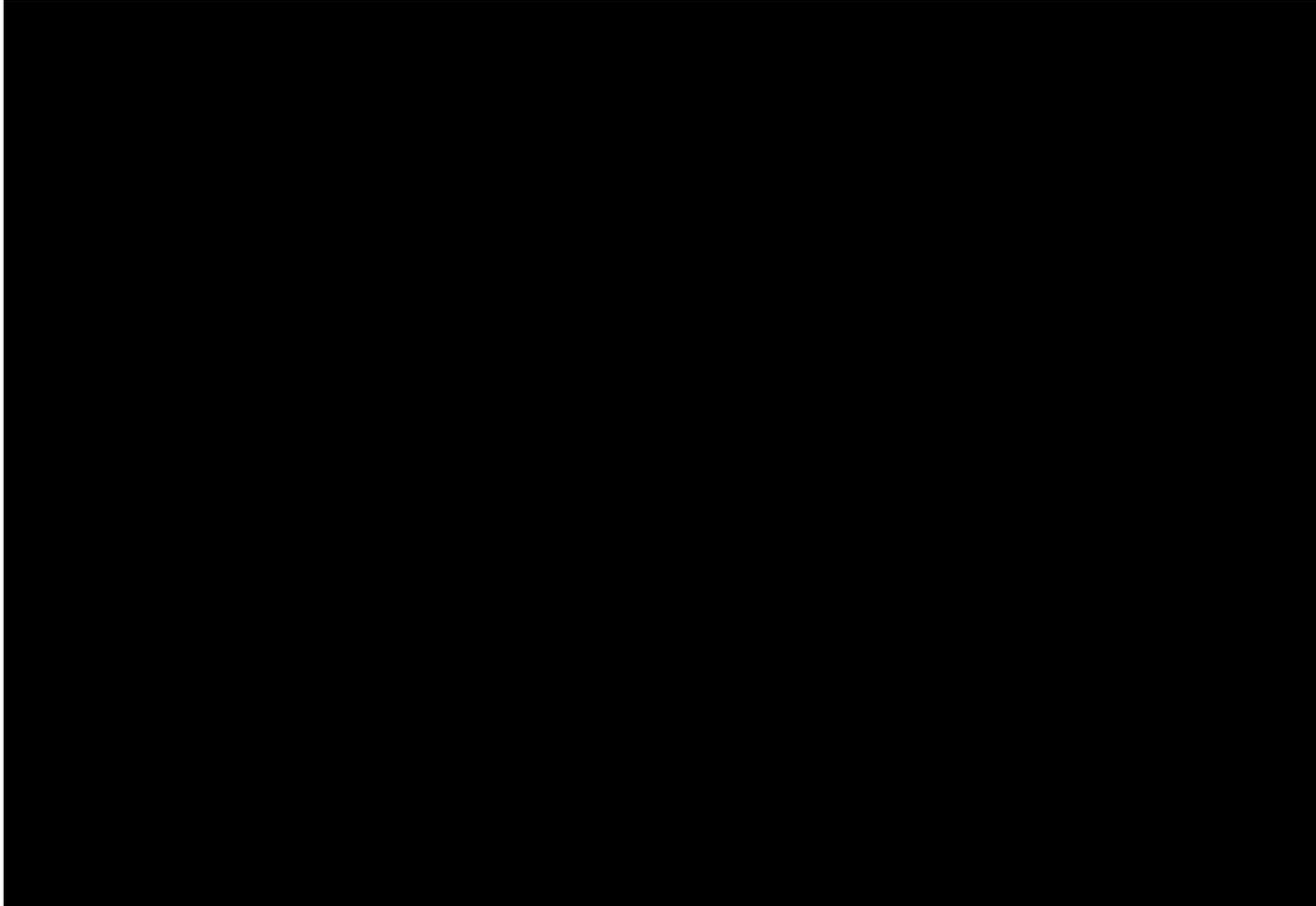
Sharing documents and teleconferencing

- Modern workflow involves web conferences
 - Google Hangout, Skype, GoToMeeting, etc.
 - People often need to **share documents**
- People share their screen, because:
 - it's **ubiquitous**: every conferencing tool has it
 - they can still **interact** while talking
 - they can share only what is **important**
 - others might not have the right software or **credentials** to open the document even if they sent it

DocuGram solution

- Let users **simply capture their screen** as they interact with the document: scrolling, mouse actions, voice comments
 - Process the recorded frames to reverse engineer the document that was shown
 - Recover **document pages**
 - Extract user's **actions** (mouse, text selections)
 - Link user's **voice** to the actions
-  Result is a sharable web link that shows the document

Demonstration



Technical Steps

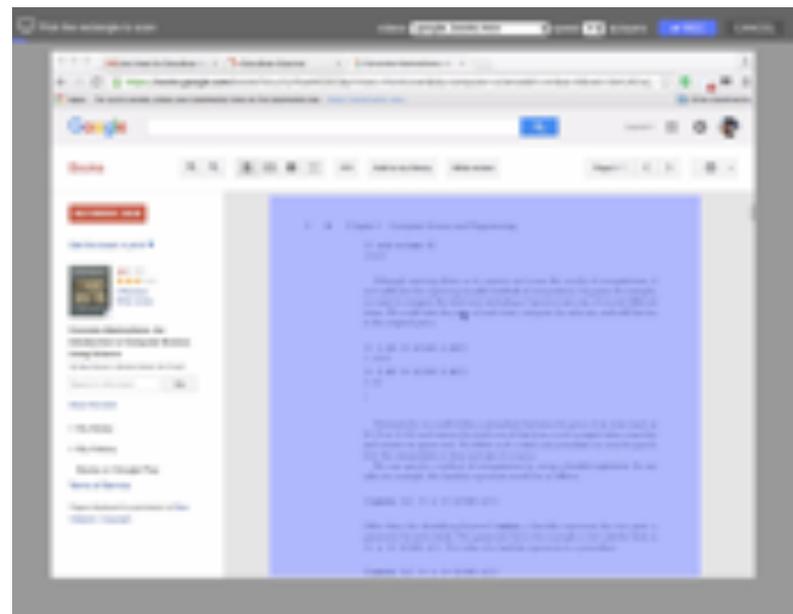
- Region of Interest (ROI) computation
- Based on image stitching
 - Fast feature point extractor
 - Robust and fast to match binary descriptors
 - Matching of detected descriptors between 2 frames
 - Gives a vertical motion ΔY
- When no vertical motion is found
 - Fast frame difference to estimate mouse position
- Stitch all frames into one image and segment tall image into pages based on horizontal breaks



Reversed engineered document from screen recording!

Region of Interest (ROI)

- Goal is to find the region corresponding to the document
- Fast binarization computes horizontal and vertical edges
- Aligned verticals are grouped, aligned horizontals are grouped
- Pairs that roughly join are candidate corners



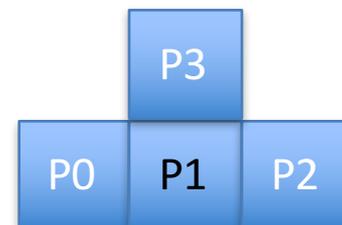
Pick the biggest of rectangles found by the corners

Image stitching (1)

- Used in video processing
 - E.g. creating panoramas in real-time
- Based on 2 insights
 - Extract consistent **key-points** from each image
 - Compute a robust **descriptor** of each keypoint
- Consistent = pick corresponding points across frames
- Robust = can tolerate some amount of noise, ~~rotation, translation~~

Image stitching (2)

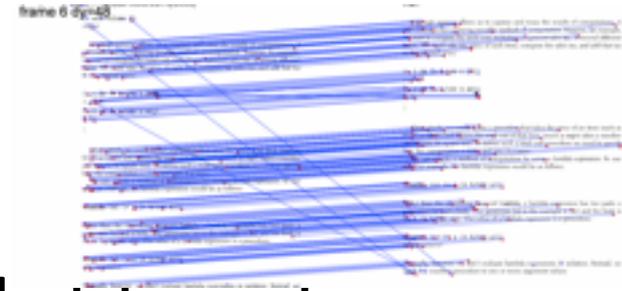
- We use a **fast** key-point extractor that works well for textual content



$$|P0-P1| > 32 \text{ AND } |P0-P2| \leq 2 \text{ AND } |P0-P3| > 32$$

- We use the BRIEF binary descriptor
 - $\text{BRIEF}_{p(x,y)} = \{0, \dots, 511\}$ bits = 64 bytes
 - $\text{Similarity}(B_{p1}, B_{p2}) = \text{XOR}(B_{p1}, B_{p2}) \leq \text{VERY FAST!!!}$

Computing vertical shift



- For every video frame captured at time t , we have key-points and descriptors:

$$(pts, desc)_{frame(t)}$$

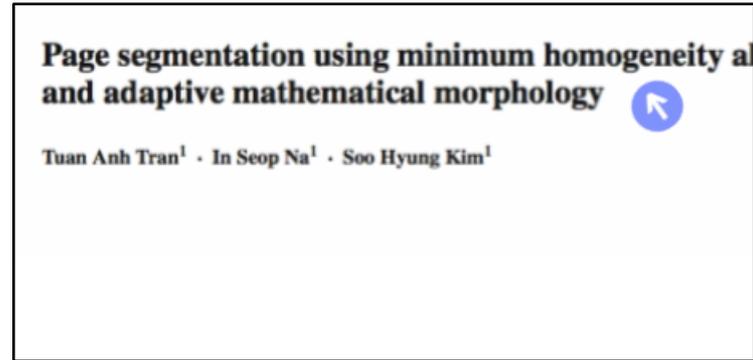
- Now loop through the pairs and find best match:

$$(pts, desc)_{frame(t-1)} \text{ with } (pts, desc)_{frame(t-1)}$$

➔ Gives us the vertical shift between the frames

Detection of mouse position

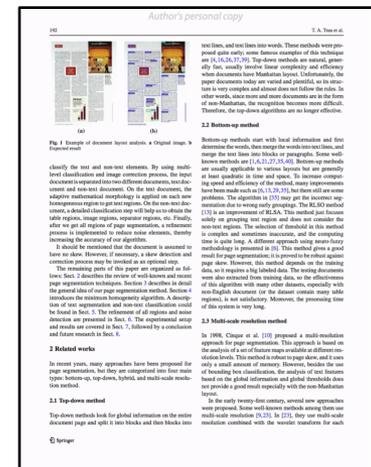
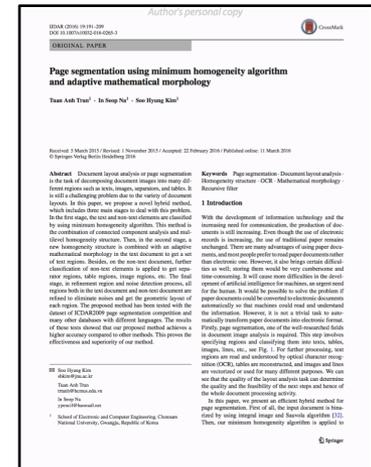
- If vertical shift is zero, compute frame difference
 - Intuitively, people don't scroll when they move the cursor
- Needs to be **fast**
 - Reduce images to 64x64 pixels before comparing



Obtain coarse location of mouse cursor

Final document re-generation

- Allocate a tall image = $\text{Sum}(\text{delta}Y)$
 - Paint each frame at the corresponding $\text{delta}Y$
 - Detect likely page boundaries
 - Binarize tall image using only vertical difference
 - Detect long horizontals and their gaps
 - Cut the long image into individual pages
- Obtains one image per page



Text detection

- Users like to select text in documents
 - XY-Cut on binarized document pages
 - Group **connected components** into text lines
 - Split each text line into word candidates based on vertical **projection profiles**

- CANVAS detects mouse motion and highlights text rectangles

 User feels interacting with the original document

Received: 5 March 2015 / Revised: 1 November 2015 / Accepted: 22 February 2016
© Springer-Verlag Berlin Heidelberg 2016

Abstract Document layout analysis or page segmentation is the task of decomposing document images into many different regions such as texts, images, separators, and tables. It is still a challenging problem due to the variety of document layouts. In this paper, we propose a novel hybrid method, which includes three main stages to deal with this problem. In the first stage, the text and non-text elements are classified by using minimum homogeneity algorithm. This method is the combination of connected component analysis and multilevel homogeneity structure. Then, in the second stage, a new homogeneity structure is combined with an adaptive mathematical morphology in the text document to get a set of text regions. Besides, on the non-text document, further classification of non-text elements is applied to get separator regions, table regions, image regions, etc. The final stage, in refinement region and noise detection process, all regions both in the text document and non-text document are

Conclusion

- **Easy-to-use tool** for users to capture and share any text-based document shown by any application
- **Real-time** processing inside their browser
- No application to install (only Chrome, Firefox)
- Novel use of **video processing** techniques to reverse engineer displayed documents
- Capture and replay of user's **interactions**

Future work

- **OCR** on captured document page images
 - Problems with low resolution, anti-aliasing
 - Promising method = **LSTM** networks for OCR
 - They don't assume pre-segmentation
- Beyond cursor motion
 - Users also select **text**, highlight passages, etc.
- Integration inside **live** video-conferencing
 - Users simply share their screen
 - DocuGram technology reverse engineers the documents and inserts them into the chat window
 - e.g. Slack, Hangout, HipChat

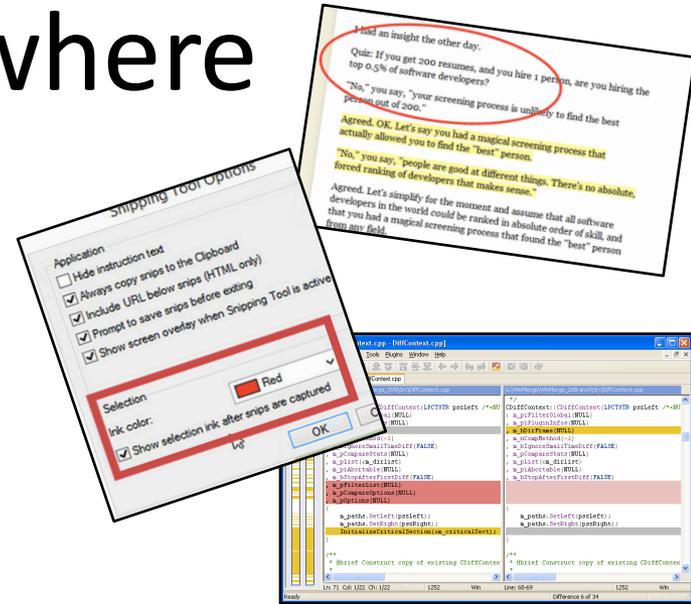
References

- Calonder, M., Lepetit, V., Strecha, C. and Fua, P., 2010. Brief: Binary robust independent elementary features. Computer Vision–ECCV 2010, pp.778-792.
- Javascript CANVAS image-manipulation and BRIEF implementation <https://trackingjs.com/>

OLD SLIDES

Screenshots are everywhere

- Everybody takes screenshots
 - Keep receipts, review web-sites



- But how do you screenshot a document?
 - Screenshot first page
 - Scroll down
 - Repeat 1 and 2
 - Stitch inside Photoshop?
- ➡ Not very easy

Problem definition

- Sharing a document requires many hoops
- It's not easy to share only a part of the document
- It's next to impossible to comment what we share
 - Export to PDF, highlight?
 - What about voice?
 - What about cursor motion?



We can find a better way

Why document screenshots?

- Digital born documents should be easy to share
 - Just email me the document
 - What format? PDF, DOC, DOCX, DOC 95, PPT, PPTX
 - Will my colleague be able to open this?
 - Can I share just 2 pages?
 - Just send the link of the document
 - Web page, Google Drive, Office 365, GitHub code
 - Will my colleague have the credentials?
 - Can I share just slides 4 and 5?

Application scenarios

- User shares a PDF document shown inside Preview App
- User shares a few slides from a PowerPoint deck
- User shares a few lines of code from his text editor

