USER INTERFACE DESIGN + PROTOTYPING + EVALUATION

**User Testing &
Automated Evaluation**

Prof. James A. Landay
University of Washington
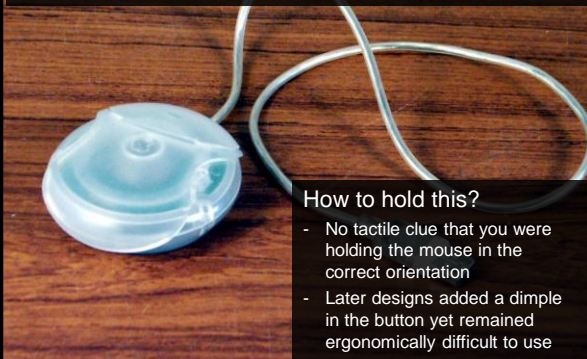
Autumn 2012

---

Product Hall of Fame or Shame?

Apple One Button Mouse

---

Product Hall of Shame!

How to hold this?
- No tactile clue that you were holding the mouse in the correct orientation
- Later designs added a dimple in the button yet remained ergonomically difficult to use

---

---

## Outline

- Visual design review
- Why do user testing?
- Choosing participants
- Designing the test
- Collecting data
- Analyzing the data
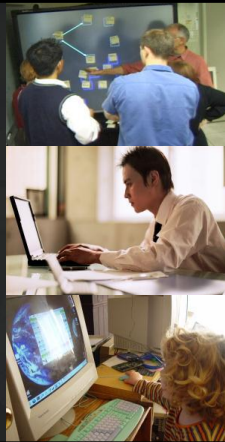- Automated evaluation

---

## Visual Design Review

- Grid systems help us put information on the page in a logical manner
  – similar things close together
- Small changes help us see key differences (e.g., small multiples)
- RGB color space leads to bad colors
- Use color properly – not for ordering!
- Avoid clutter – remove until you can remove no more

## Why do User Testing?

- Can't tell how good UI is until?
  - people use it!

- Expert review methods are based on evaluators who?
  - may know too much
  - may not know enough (about tasks, etc.)

- Hard to predict what real users will do

## Choosing Participants

- Representative of target users
  - job-specific vocab / knowledge
  - tasks
- Approximate if needed
  - system intended for doctors?
    - get medical students or nurses
  - system intended for engineers?
    - get engineering students
- Use incentives to get participants

## Ethical Considerations

- Usability tests can be distressing
  - users have left in tears

- You have a responsibility to alleviate
  - make voluntary with informed consent (form)
  - avoid pressure to participate
  - let them know they can stop at any time
  - stress that you are testing the system, not them
  - make collected data as anonymous as possible
- Often must get human subjects approval

11/27/2012 CSE 440: User Interface Design, Prototyping, & Evaluation 9

## User Test Proposal

- A report that contains
  - objective
  - description of system being testing
  - task environment & materials
  - participants
  - methodology
  - tasks
  - test measures

- Get approved & then reuse for final report
- Seems tedious, but writing this will help "debug" your test

11/27/2012 CSE 440: User Interface Design, Prototyping, & Evaluation 10

## Selecting Tasks

- Should reflect what real tasks will be like
- Tasks from analysis & design can be used
  - may need to shorten if
    - they take too long
    - require background that test user won't have

- Try not to train unless that will happen in real deployment
- Avoid bending tasks in direction of what your design best supports
- Don't choose tasks that are too fragmented
  - e.g., phone-in bank test

11/27/2012 CSE 440: User Interface Design, Prototyping, & Evaluation 11

## Two Types of Data to Collect

- Process data
  - observations of what users are doing & thinking

- Bottom-line data
  - summary of what happened (time, errors, success)
  - i.e., the dependent variables

11/27/2012 CSE 440: User Interface Design, Prototyping, & Evaluation 12

## Which Type of Data to Collect?

- Focus on process data first
  - gives good overview of where problems are
- Bottom-line doesn't tell you ?
  - where to fix
  - just says: "too slow", "too many errors", etc.
- Hard to get reliable bottom-line results
  - need many users for statistical significance

## The "Thinking Aloud" Method

- Need to know what users are thinking, not just what they are doing
- Ask users to talk while performing tasks
  - tell us what they are thinking
  - tell us what they are trying to do
  - tell us questions that arise as they work
  - tell us things they read
- Make a recording or take good notes
  - make sure you can tell what they were doing

## Thinking Aloud (cont.)

- Prompt the user to keep talking
  - "tell me what you are thinking"

- Only help on things you have pre-decided
  - keep track of anything you do give help on

- Recording
  - use a digital watch/clock
  - take notes, plus if possible
    - record audio & video (or even event logs)

## Using the Test Results

- Summarize the data
  - make a list of all critical incidents (CI)
    - positive & negative
  - include references back to original data
  - try to judge why each difficulty occurred

- What does data tell you?
  - UI work the way you thought it would?
    - users take approaches you expected?
  - something missing?

## Using the Results (cont.)

- Update task analysis & rethink design
  - rate severity & ease of fixing CIs
  - fix both severe problems & make the easy fixes

## Will thinking out loud give the right Answers?

- Not always

- If you ask a question, people will always give an answer, even it is has nothing to do with facts
  - panty hose example

→Try to avoid specific questions

## Analyzing the Numbers

- Example: trying to get task time ≤ 30 min.
  - test gives: 20, 15, 40, 90, 10, 5
  - mean (average) = 30
  - median (middle) = 17.5
  - looks good!
- Did we achieve our goal?
- Wrong answer, not certain of anything!
- Factors contributing to our uncertainty?
  - small number of test users (n = 6)
  - results are very variable (standard deviation = 32)
    - std. dev. measures dispersal from the mean

## Measuring Bottom-Line Usability

- Situations in which numbers are useful
  - time requirements for task completion
  - successful task completion %
  - compare two designs on speed or # of errors
- Ease of measurement
  - time is easy to record
  - error or successful completion is harder
    - define in advance what these mean
- Do not combine with thinking-aloud. Why?
  - talking can affect speed & accuracy

## Analyzing the Numbers (cont.)

- This is what statistics is for

- Crank through the procedures and you find
  - 95% certain that typical value is between 5 & 55

## Analyzing the Numbers (cont.)

| Web Usability Test Results | | | | |
|---|---|---|---|---|
| Participant # | Time (minutes) | | | |
| 1 | 20 | | | |
| 2 | 15 | | | |
| 3 | 40 | | | |
| 4 | 90 | | | |
| 5 | 10 | | | |
| 6 | 5 | | | |
| | | | | |
| number of participants | 6 | | | |
| mean | 30.0 | | | |
| median | 17.5 | | | |
| std dev | 31.8 | | | |
| | | | | |
| standard error of the mean | = stddev / sqrt (#samples) | 13.0 | | |
| | | | | |
| typical values will be mean +/- 2*standard error | --> 4 to 56! | | | |
| | | | | |
| what is plausible? = confidence (alpha=5%, stddev, sample size) | 25.4 | --> 95% confident between 5 & 56 | | |

## Analyzing the Numbers (cont.)

- This is what statistics is for

- Crank through the procedures and you find
  - 95% certain that typical value is between 5 & 55

- Usability test data is quite variable
  - need lots to get good estimates of typical values
  - 4 times as many tests will only narrow range by 2x
    - breadth of range depends on sqrt of # of test users
  - this is when online methods become useful
    - easy to test w/ large numbers of users

## Measuring User Preference

- How much users like or dislike the system
  - can ask them to rate on a scale of 1 to 10
  - or have them choose among statements
    - "best UI I've ever…", "better than average"…
  - hard to be sure what data will mean
    - novelty of UI, feelings, not realistic setting …
- If many give you low ratings → trouble

- Can get some useful data by asking
  - what they liked, disliked, where they had trouble, best part, worst part, etc.
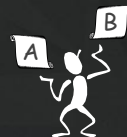  - redundant questions are OK

11/27/2012    CSE 440: User Interface Design, Prototyping, & Evaluation    25

## Comparing Two Alternatives

- *Between groups* experiment

  - two groups of test users
  - each group uses only 1 of the systems

- *Within groups* experiment

  - one group of test users
    - each person uses both systems
    - can't use the same tasks or order (learning)
  - best for low-level interaction techniques

11/27/2012    CSE 440: User Interface Design, Prototyping, & Evaluation    26

## Comparing Two Alternatives

- Between groups requires many more participants than within groups

- See if differences are statistically significant
  - assumes normal distribution & same std. dev.

- Online companies can do large AB tests
  - look at resulting behavior (e.g., buy?)

11/27/2012    CSE 440: User Interface Design, Prototyping, & Evaluation    27

## Experimental Details

- Order of tasks
  - choose one simple order (simple → complex)
    - unless doing within groups experiment

- Training
  - depends on how real system will be used

- What if someone doesn't finish
  - assign very large time & large # of errors or remove & note

- Pilot study
  - helps you fix problems with the study
  - do two, first with colleagues, then with real users

11/27/2012    CSE 440: User Interface Design, Prototyping, & Evaluation    28

## Instructions to Participants

- Describe the purpose of the evaluation
  - "I'm testing the product; I'm not testing you"
- Tell them they can quit at any time
- Demonstrate the equipment
- Explain how to think aloud
- Explain that you will not provide help
- Describe the task

  - give written instructions, one task at a time

11/27/2012    CSE 440: User Interface Design, Prototyping, & Evaluation    29

## Details (cont.)

- Keeping variability down
  - recruit test users with similar background
  - brief users to bring them to common level
  - perform the test the same way every time
    - don't help some more than others (plan in advance)
  - make instructions clear

- Debriefing test users
  - often don't remember, so demonstrate or show video segments
  - ask for comments on specific features
    - show them screen (online or on paper)

11/27/2012    CSE 440: User Interface Design, Prototyping, & Evaluation    30

## Reporting the Results

- Report what you did & what happened
- Images & graphs help people get it!
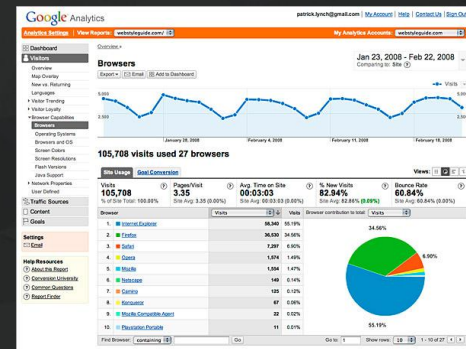- Video clips can be quite convincing



## AUTOMATED & REMOTE USABILITY EVALUATION

## Automated Analysis & Remote Testing

- Log analysis
  - infer user behavior by looking at web server logs

- A-B Testing
  - show different user segments different designs
  - requires live site (built) & customer base
  - measure outcomes (profit), but not why?

- Remote user testing
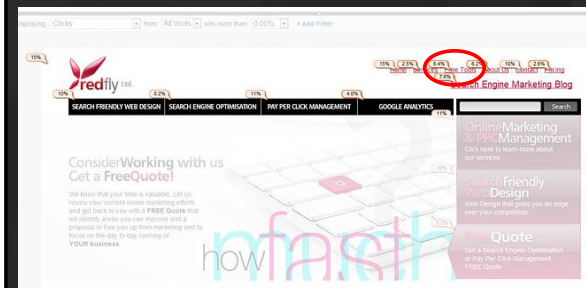  - similar to in lab, but online (e.g., over Skype)

11/27/2012     CSE 440: User Interface Design, Prototyping, & Evaluation     33

## Web Logs Analysis Difficult



11/27/2012     CSE 440: User Interface Design, Prototyping, & Evaluation     34

## Google Analytics – Server Logs++



http://www.redflymarketing.com/blog/using-google-analytics-to-improve-conversions/

11/27/2012     CSE 440: User Interface Design, Prototyping, & Evaluation     35

## Google Analytics – Server Logs++



http://www.redflymarketing.com/blog/using-google-analytics-to-improve-conversions/

11/27/2012     CSE 440: User Interface Design, Prototyping, & Evaluation     36

THE WALL STREET JOURNAL.
WSJ.com

AUGUST 17, 2009
INNOVATION

The New, Faster Face of Innovation

*Thanks to technology, change has never been so easy—or so cheap*

By ERIK BRYNJOLFSSON And MICHAEL SCHRAGE

Call it innovation on steroids. Or innovation at warp speed. Or just the innovation of rapid innovation.

But the essential point remains: Technology is transforming innovation at its core, allowing companies to test new ideas at speeds—and prices—that were unimaginable even a decade ago. They can stick features on Web sites and tell within hours how customers respond. They can see results from in-store promotions, or efforts to boost process productivity, almost as quickly.

## Web Allows Controlled A/B Experiments



- Example: Amazon Shopping Cart
  - Add item to cart
  - Site shows cart contents
- Idea: show recommendations based on cart items
- Arguments
  - Pro: cross-sell more items
  - Con: distract people at check out
- Highest Paid Person's Opinion *"Stop the project!"*
- Simple experiment was run, wildly successful

From Greg Linden's Blog: http://glinden.blogspot.com/2006/04/early-amazon-shopping-cart.html

11/27/2012 — CSE 440: User Interface Design, Prototyping, & Evaluation — 38

## Windows Marketplace: *Solitaire vs. Poker*

**Which image has the higher clickthrough? By how much?**

A: Solitaire game

B: Poker game



**A is 61% better. Why?**

Courtesy of Ronny Kohavi

## The Trouble With Most Web Site Analysis Tools



*Leave*

Unknowns
- Who?
- What?
- Why?
- Did they find it?
- Satisfied?

11/27/2012 — CSE 440: User Interface Design, Prototyping, & Evaluation — 40

## NetRaker Usability Research
### See how customers accomplish real tasks on site



## NetRaker Usability Research
### See how customers accomplish real tasks on site

## NetRaker Usability Research
### See how customers accomplish real tasks on site



## UserZoom



---

## Advantages of Remote Usability Testing

- Fast
  - can set up research in 3-4 hours
  - get results in 36 hours
- More accurate
  - can run with large samples (50-200 users → stat. sig.)
  - uses real people (customers) performing tasks
  - natural environment (home/work/machine)
- Easy-to-use
  - templates make setting up easy
- Can compare with competitors
  - indexed to national norms

---

## Disadvantages of Remote Usability Testing

- Miss observational feedback
  - facial expressions
  - verbal feedback (critical incidents)

- Need to involve human participants
  - costs some amount of money (typically $20-$50/person)

- People often do not like pop-ups
  - need to be careful when using them

---

## Summary

- User testing is important, but takes time/effort
- Early testing can be done on mock-ups (low-fi)
- Use ????? tasks & ????? participants
  - real tasks & representative participants
- Be ethical & treat your participants well
- Want to know what people are doing & why? collect
  - process data
- Bottom line data requires ???? to get statistically reliable results
  - more participants
- Difference between between & within groups?
  - between groups: everyone participates in one condition
  - within groups: everyone participates in multiple conditions
- Automated usability
  - faster than traditional techniques
  - can involve more participants → convincing data
  - easier to do comparisons across sites
  - tradeoff with losing observational data

---

## Next Time

Interactive Prototype Presentations