## Evaluating Interfaces with Users

**Why evaluation is crucial to interface design**

**General approaches and tradeoffs in evaluation**
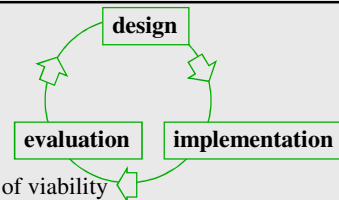
**The role of ethics**

---

## Why Bother?

**design**

**evaluation**　　**implementation**

**Tied to the usability engineering lifecycle**

- Pre-design
  - investing in new expensive systems requires proof of viability

- Initial design stages
  - develop and evaluate initial design ideas with the user

- Iterative design
  - does system behaviour match the user's task requirements?
  - are there specific problems with the design?
  - can users provide feedback to modify design?

- Acceptance testing
  - verify that human/computer system meets expected performance criteria
    - ease of learning, usability, user's attitude, performance criteria
    - e.g., a first time user will take 1-3 minutes to learn how to withdraw $50 from the automatic teller

Evan Golub / Ben Bederson / Saul Greenberg

### What Defines Success?

**We want a "usable" system. What are some metrics that can be used to measure whether a system is usable?**
- Time to learn
- Speed of performance
- Rate of errors by users
- Retention over time
- Subjective Satisfaction

**Often, there will be tradeoffs between these goals.**

### Approaches: Naturalistic/Qualitative

**Naturalistic:**
- describes an ongoing process as it evolves over time
- observation occurs in realistic setting
  - ecologically valid
- "real life"

**External validity**
- degree to which research results applies to real situations

## Approaches: Experimental/Quantitative

**Experimental**
- study relations by manipulating one or more *independent* variables
  - experimenter controls all environmental factors
- observe effect on one or more *dependent* variables

**Internal validity**
- confidence that we have in our explanation of experimental results

**Trade-off: Natural *vs* Experimental**
- precision and direct control over experimental design
  *versus*
- desire for maximum generalizability in real life situations

## Reliability Concerns

**Would the same results be achieved if the test were repeated?**

**Problem: individual differences:**
- best user 10x faster than slowest
- best 25% of users ~2x faster than slowest 25%

**Partial Solution**
- reasonable number and range of users tested
- statistics provide confidence intervals of test results
  - 95% confident that mean time to perform task X is 4.5+/-0.2 minutes means
    95% chance true mean is between 4.3 and 4.7, 5% chance its outside that

## Validity Concerns

**Does the test measure something of relevance to usability of real products in real use outside of lab?**

- Some typical validity problems of testing vs real use
  – non-typical users tested
  – tasks are not typical tasks
  – physical environment different
      quiet lab -vs- very noisy open offices vs interruptions
  – social influences different
      motivation towards experimenter vs motivation towards boss

**Partial Solution**

- use real users
- tasks from task-centered system design
- environment similar to real situation

## Qualitative methods for usability evaluation

**Qualitative:**

- produces a description, usually in non-numeric terms
- may be subjective

**Methods**

- **Introspection**
  – by designer
  – by users
- **Direct observation**
  – simple observation
  – think-aloud
  – constructive interaction
- **Query**
  – interviews (structured and retrospective)
  – surveys and questionnaires

## Introspection Method

## Introspection Method: Designer

**The designer tries the system (or prototype) out (a walkthrough of the systems screens and features)**

- does the system "feel right"?
- most common evaluation method

**Problems**
  – not reliable as completely subjective
  – not valid as "introspector" is a non-typical user

**Intuitions and introspection are often wrong!**

## Introspection Method: User

**Conceptual Model Extraction**

- Show the users low-fidelity prototypes or screenshots of medium-fidelity prototypes (user-centered walkthrough).
- Ask the user to explain what each screen element does or represents as well as how they would attempt to perform individual tasks.
- This allows us to gain insight as to a user's initial perception of our interface and the mental model they might be constructing as they begin to use our system.

**NOTE: Since we are walking them through specific parts as their guide, we will not really see how a user might explore the system on their own or their learning processes.**

## Direct observation (three approaches within here)

**Evaluator observes and records users interacting with design/system**

- in lab:
  - user asked to complete a set of pre-determined tasks
  - a specially built and fully instrumented usability lab may be available
- in field:
  - user goes through normal duties

**Excellent at identifying gross design/interface problems**

**Validity/reliability depends on how controlled/contrived the situation is...**

**Three general approaches:**
- simple observation
- think-aloud
- constructive interaction

## Direct observation: Simple Observation Method

**User is given the task, and evaluator just watches the user**

**Problem**

- does not give insight into the user's decision process or attitude

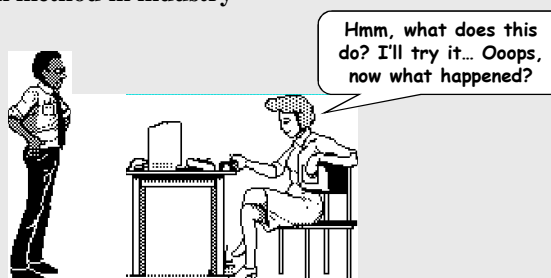## Direct observation: Think Aloud Method

**Subjects are asked to say what they are thinking/doing**
- – what they believe is happening
- – what they are trying to do
- – why they took an action
- Gives insight into what the user is thinking

**Problems**
- – awkward/uncomfortable for subject (thinking aloud is not normal!)
- – "thinking" about it may alter the way people perform their task
- – hard to talk when they are concentrating on problem

**Most widely used evaluation method in industry**

Hmm, what does this do? I'll try it… Ooops, now what happened?

## Direct observation: Constructive Interaction Method

**Two people work together on a task**

- normal conversation between the two users is monitored
  – removes awkwardness of think-aloud
- Variant: Co-discovery learning
  – use semi-knowledgeable "coach" and naive subject together
  – make naive subject use the interface
- results in
  – naive subject asking questions
  – semi-knowledgeable coach responding
  – provides insights into thinking process of both beginner and intermediate users



> Now, why did it do that?

> Oh, I think you clicked on the wrong icon

---

## Recording Observations

**Make sure you get permission!**

**Make sure you are mindful of privacy!**

## **Recording Observations**: Tools

**How do we record user actions during observation for later analysis?**
  – if no record is kept, evaluator may forget, miss, or misinterpret events

- paper and pencil
  – primitive but cheap
  – evaluators record events, interpretations, and extraneous observations
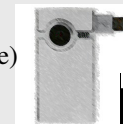  – hard to get detail (writing is slow)
  – coding schemes help…
- audio recording
  – good for recording talk produced by thinking aloud/constructive interaction
  – hard to tie into user actions (ie what they are doing on the screen)
  – hard to search through later
- video recording
  – can see and hear what a user is doing
  – one camera for screen, another for subject (picture in picture)
  – can be intrusive during initial period of use
  – generates too much data

---

## **Example coding scheme...**

**Tracking a person's activity in the office with quick notations.**

s = start of activity
e = end of activity

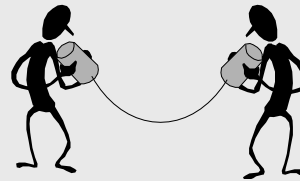| Time | Desktop Activities | | | Absences from Desk | | Interruptions | | |
|------|----------------------|------------------|-------------------|---------|-------------|--------|-------|--------|
|      | Working on computer | Working at desk | Using telephone | In room | Out of room | Person | Phone | e-mail |
| 9:00 | s |   |   |   |   |   |   |   |
| 9:02 | e |   |   |   |   | s |   |   |
| 9:05 |   |   |   |   |   | e |   |   |
| 9:10 |   |   |   |   | s |   |   |   |
| 9:13 |   |   | s |   | e |   |   |   |
|      |   |   |   |   |   |   |   |   |

## Querying Users: Interviews

**Excellent for pursuing specific issues**
- vary questions to suit the context
- probe more deeply on interesting issues as they arise
- good for exploratory studies via open-ended questioning
- often leads to specific constructive suggestions

**Problems:**
- accounts are subjective
- time consuming
- evaluator can easily bias the interview
- prone to rationalization of events/thoughts by user
  – user's reconstruction may be wrong

## Querying Users: Structured Interviews

**Plan a set of central questions**
- could be based on results of user observations
- gets things started
- focuses the interview
- ensures a base of consistency

**Try not to ask leading questions!**
"Now that was easy, wasn't it?"
"How hard would you say this task was?"

**Start with individual discussions to discover different perspectives,
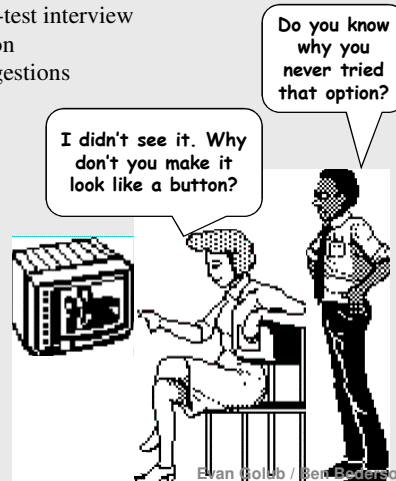and continue with group discussions**
- the larger the group, the more the universality of comments can be ascertained
- also encourages discussion between users

## Querying Users: Retrospective Testing

**Post-observation interview to clarify events that occurred during system use**
- perform an observational test
- create a video record of it
- have users view the video and comment on what they did
    - excellent for grounding a post-test interview
    - avoids erroneous reconstruction
    - users often offer concrete suggestions

Do you know why you never tried that option?

I didn't see it. Why don't you make it look like a button?

---

## Querying Users: Surveys and Questionnaires

**Preparation "expensive," but administration cheap**
- can reach a wide subject group (e.g. mail)

**Does not require presence of evaluator.**

**Results can be quantified.**

**Only as good as the questions asked!!!**

**Often has low return rate (what's in it for them?) or biased sample (who will take the time to answer?)**

QUIS - Questionnaire for User Interface Satisfaction
- Shneiderman/Plaisant text has some example questions from it on pages 152-154.

## Querying Users: Surveys and Questionnaires Details

**Establish the _purpose_ of the questionnaire**
- what information is sought?
- how would you analyze the results?
- what would you do with your analysis?

**Typically will not ask questions whose answers you will not use**
- this is unlike many other types of surveys you may have discussed in your psychology class

**Determine the _audience_ you want to reach**
- typical survey: random sample of between 50 and 1000 users of the product

**Determine how would you will deliver and collect the questionnaire**
- on-line for computer users
- surface mail (with pre-addressed reply envelope for better response rate)

**Determine target demographics**
- e.g. level of experience, age, income, etc.

## Styles of Questions (I)

**Open-ended questions**
- asks for unprompted opinions
- good for general subjective information
  – but difficult to analyze rigorously

eg: **Can you suggest any improvements to the interfaces?**

## Styles of Questions (II)

**Closed questions**
- restricts the respondent's responses by supplying alternative answers
- makes questionnaires a chore for respondent to fill in
- can be easily analyzed
- but watch out for hard to interpret responses!
  – alternative answers should be very specific

  Do you use computers at work:
    O often           O sometimes     O rarely
                    *-vs-*
  In your typical work day,  do you use computers:
    O over 4 hrs a day
    O between 2 and 4 hrs daily
    O between 1and 2 hrs daily
    O less than 1 hr a day

## Styles of Questions (III)

**Bipolar Scaling**
- ask user to judge a specific statement on a numeric scale
- scale usually corresponds with agreement or disagreement with a statement

  Characters on the computer screen are:
    hard to read  **1**  **2**  **③**  **4**  **5**  easy to read

Scale of **1 to 7** or **1 to 9** might provide better results since they will still
  provide a good range even if the user eliminates the extremes.

Sometimes done explicitly as:
1. Strongly disagree
2. Disagree
3. Neutral
4. Agree
5. Strongly agree

Scale which is **even** in length should be used if you want to prevent the user
  from being neutral.

## Styles of Questions (IV)

**Multiple choice (possibly multiple responses)**
  • respondent offered a choice of explicit responses

How do you most often get help with the system? (tick one)
O    on-line manual
O    paper manual
Ø    ask a colleague

Which types of software have you used? (tick all that apply)
O    word processor
Ø    data base
O    spreadsheet
Ø    compiler

## Styles of Questions (V)

**Ranked**
  • respondent places an ordering on items in a list
  • useful to indicate a user's preferences
  • forced choice

Rank the usefulness of these methods of issuing a command
(1 most useful, 2 next most useful..., 0 if not used
__2__ command line
__1__ menu selection
__3__ control key accelerator

## Styles of Questions (VI)

**Combining open-ended and closed questions**

• gets specific response, but allows room for user's opinion

It is easy to recover from mistakes:

disagree                     agree        comment: *the undo facility is really*
  *helpful*
      1     2     3    ④    5

## What might the future hold?

We live in a time where the use of AI is on the rise and chat bots are in the thick of things.  We've gone from Universal McCann's **Jill020306** to Microsoft's **Tay** in the span of a decade.  Where could they take us in the future?

When you write a new library or program module, you can use unit testing tools to automatically assess the accuracy of various things.  The tools continue to expand, even into the ability to automatically test GUI elements.  Where might they go in the future?

**Possible direction: "Chatbots: Your Ultimate Prototyping Tool"**

https://medium.com/ideo-stories/chatbots-ultimate-prototyping-tool-e4e2831967f3#.74ci9jy2n

## What you now know about…

**Observing a range of users use your system for specific tasks can reveal successes and problems and qualitative observational tests can be quick (and somewhat easy) to do. Several methods can reveal what is in a person's head as they are doing the test. Particular methods include:**

- Conceptual model extraction
- Direct observation (simple observation, think-aloud, constructive interaction)
- Query via interviews, retrospective testing and questionnaires
- Continuous evaluation via user feedback and field studies

**Evaluation is crucial for designing, debugging, and verifying interfaces**

**There is a tradeoff in naturalistic *-vs-* experimental approaches**
- internal and external validity
- reliability
- precision
- generalizability

### UP NEXT: ETHICS!

## Readings

Optional reading is "Designing the User Interface" Chapter 5