



# Useable Security and Privacy

## Usability & Evaluation

Thanks go to:

- Prof. Dr. Michael Rohs



# USABILITY



# User – Tool – Task/Goal – Context



user



context



tool



task/goals



# Usability (ISO 9241 Standard)

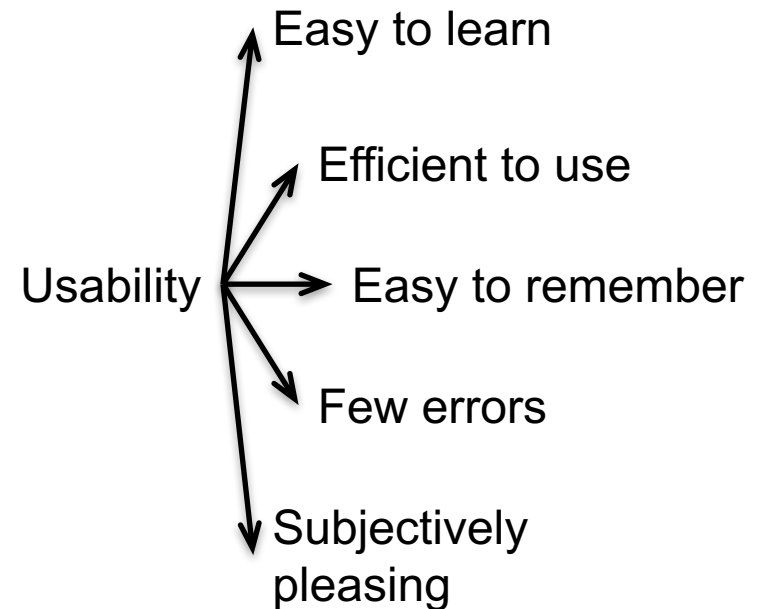
- Extent to which a product can be used by specified users to achieve specified goals with effectiveness, efficiency, and satisfaction in a specified context of use.
  - Effectiveness: Quality, accuracy, and completeness with which users achieve goals
  - Efficiency: Effort necessary to reach a certain level of quality, accuracy, and completeness
  - Satisfaction: Comfort and acceptability of the system to its users (enjoyable, motivating? or limiting, irritating?)
  - Context of use: Users, tasks, equipment, physical and social environment, organizational requirements

ISO 9241-11. Ergonomic requirements for office work with visual display terminals (VDTs)-Part 11: Guidance on usability—Part 11 (ISO 9241-11:1998)



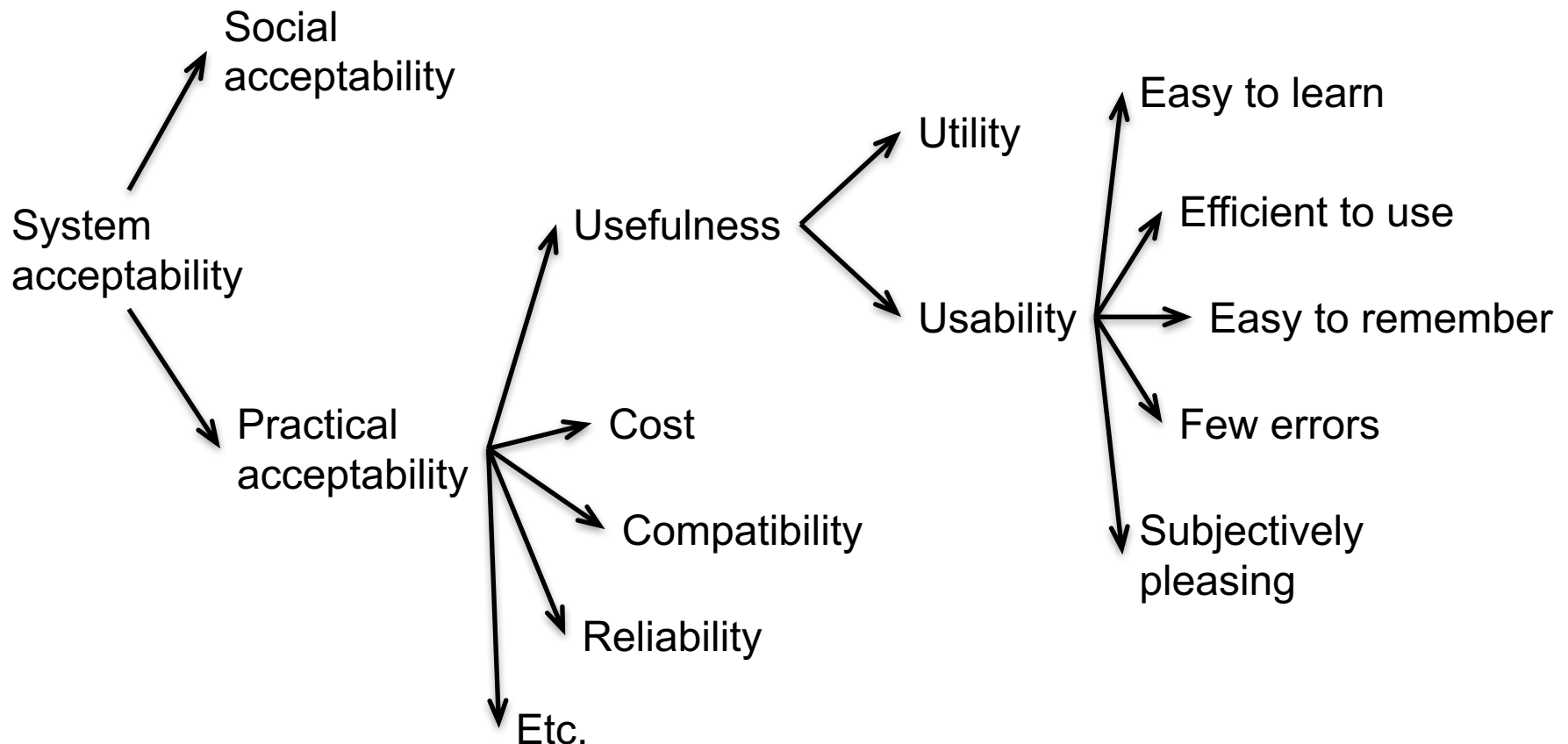
# Attributes of Usability (Nielsen)

- Learnability (easy to learn)
- Efficiency (efficient to use)
- Memorability (easy to remember)
- Errors (few errors)
- Satisfaction (subjectively pleasing)





# Usability as an Aspect of System Acceptability (Nielsen)





# Typical Measures of Effectiveness

- Binary task completion
- Accuracy
  - Error rates
  - Spatial accuracy
  - Precision
- Recall
- Completeness
- Quality of outcome
  - Understanding
  - Experts' assessment
  - Users' assessment

Kasper Hornbæk: Current practice in measuring usability: Challenges to usability studies and research. *Int. J. Human-Computer Studies* 64 (2006) 79–102.



- Time
  - Task completion time
  - Time in mode (e.g., time in help)
  - Time until event (e.g., time to react to warning)
- Input rate (e.g., words per minute, WPM)
- Mental effort (NASA Task Load Index)
  - <http://www.keithv.com/software/nasatlx/>
- Usage patterns
  - Use frequency (e.g., number of button clicks)
  - Information accessed (e.g., number of Web pages visited)
  - Deviation from optimal solution (e.g. path length)
- Learning (e.g., shorter task time over sessions)

Kasper Hornbæk: Current practice in measuring usability: Challenges to usability studies and research. *Int. J. Human-Computer Studies* 64 (2006) 79–102.





- Standard questionnaires (e.g., SUS, QUIS, AttrakDiff)
- Preference
  - Rate or rank interfaces
  - Behavior in interaction (e.g., observe what users choose)
- Satisfaction with the interface
  - Ease-of-use (e.g. 5-/7-point Likert scale: “X was easy to use”)
  - Satisfaction with specific features
  - During use (e.g., heart period variability, reflex responses)
- Attitudes and perceptions
  - Attitudes towards others (e.g., “I felt connected to X when using...”)
  - Perception of outcome / interaction

Kasper Hornbæk: Current practice in measuring usability: Challenges to usability studies and research. *Int. J. Human-Computer Studies* 64 (2006) 79–102.



# Typical Measures of Specific Attitudes

- Annoyance
- Anxiety
- Complexity
- Control
- Engagement
- Flexibility
- Fun
- Liking
- Wanting to use again

Kasper Hornbæk: Current practice in measuring usability: Challenges to usability studies and research. *Int. J. Human-Computer Studies* 64 (2006) 79–102.



# SUS: System Usability Scale

- Developed by DEC Corporation
- 10 5-point Likert scales
- Single score (0-100)
  - Odd items: position – 1
  - Even items: 5 – position
  - Add item scores
  - Multiply by 2.5

	Strongly disagree					Strongly agree
1. I think that I would like to use this system frequently	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
	1	2	3	4	5	
2. I found the system unnecessarily complex	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
	1	2	3	4	5	
3. I thought the system was easy to use	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
	1	2	3	4	5	
4. I think that I would need the support of a technical person to be able to use this system	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
	1	2	3	4	5	
5. I found the various functions in this system were well integrated	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
	1	2	3	4	5	
6. I thought there was too much inconsistency in this system	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
	1	2	3	4	5	
7. I would imagine that most people would learn to use this system very quickly	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
	1	2	3	4	5	
8. I found the system very cumbersome to use	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
	1	2	3	4	5	
9. I felt very confident using the system	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
	1	2	3	4	5	
10. I needed to learn a lot of things before I could get going with this system	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
	1	2	3	4	5	

Brooke. SUS: A "quick and dirty" usability scale. Usability Evaluation in Industry. London: Taylor and Francis, 1996



# Example: SUS-Ratings

	Strongly disagree				Strongly agree	
1. I think that I would like to use this system frequently		<b>X</b>				<b>1</b>
	1	2	3	4	5	
2. I found the system unnecessarily complex	<b>X</b>					<b>4</b>
	1	2	3	4	5	
3. I thought the system was easy to use		<b>X</b>				<b>1</b>
	1	2	3	4	5	
4. I think that I would need the support of a technical person to be able to use this system			<b>X</b>			<b>2</b>
	1	2	3	4	5	
5. I found the various functions in this system were well integrated	<b>X</b>					<b>0</b>
	1	2	3	4	5	

pos=2: score = pos-1=1

pos=1: score = 5-pos=4

pos=2: score = pos-1=1

pos=3: score = 5-pos=2

pos=1: score = pos-1=0

Brooke. SUS: A "quick and dirty" usability scale. Usability Evaluation in Industry. London: Taylor and Francis, 1996



# Example: SUS-Ratings

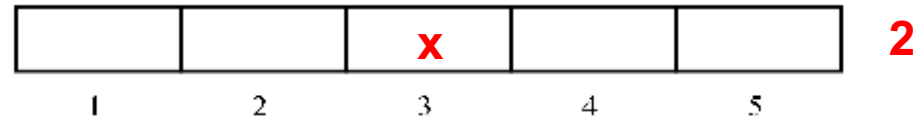
6. I thought there was too much inconsistency in this system



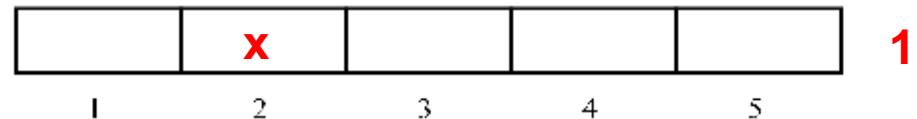
7. I would imagine that most people would learn to use this system very quickly



8. I found the system very cumbersome to use



9. I felt very confident using the system



10. I needed to learn a lot of things before I could get going with this system



**Sum = 16**

**SUS-Score = Sum \* 2.5 = 40**

Brooke. SUS: A "quick and dirty" usability scale. Usability Evaluation in Industry. London: Taylor and Francis, 1996



# QUIS: Questionnaire for User Interaction Satisfaction

- Developed by the University of Maryland
  - Academic License \$200, Student License \$50
- Semantic differential scales
- Components: (1) demographics, (2) overall reaction ratings (6 scales), (3) specific interface factors: screen, terminology and system feedback, learning, system capabilities, (4) optional sections
- Long and short forms
- <http://lap.umd.edu/quis/>

frustrating		satisfying							
1	2	3	4	5	6	7	8	9	NA
dull		stimulating							
1	2	3	4	5	6	7	8	9	NA
difficult		easy							
1	2	3	4	5	6	7	8	9	NA

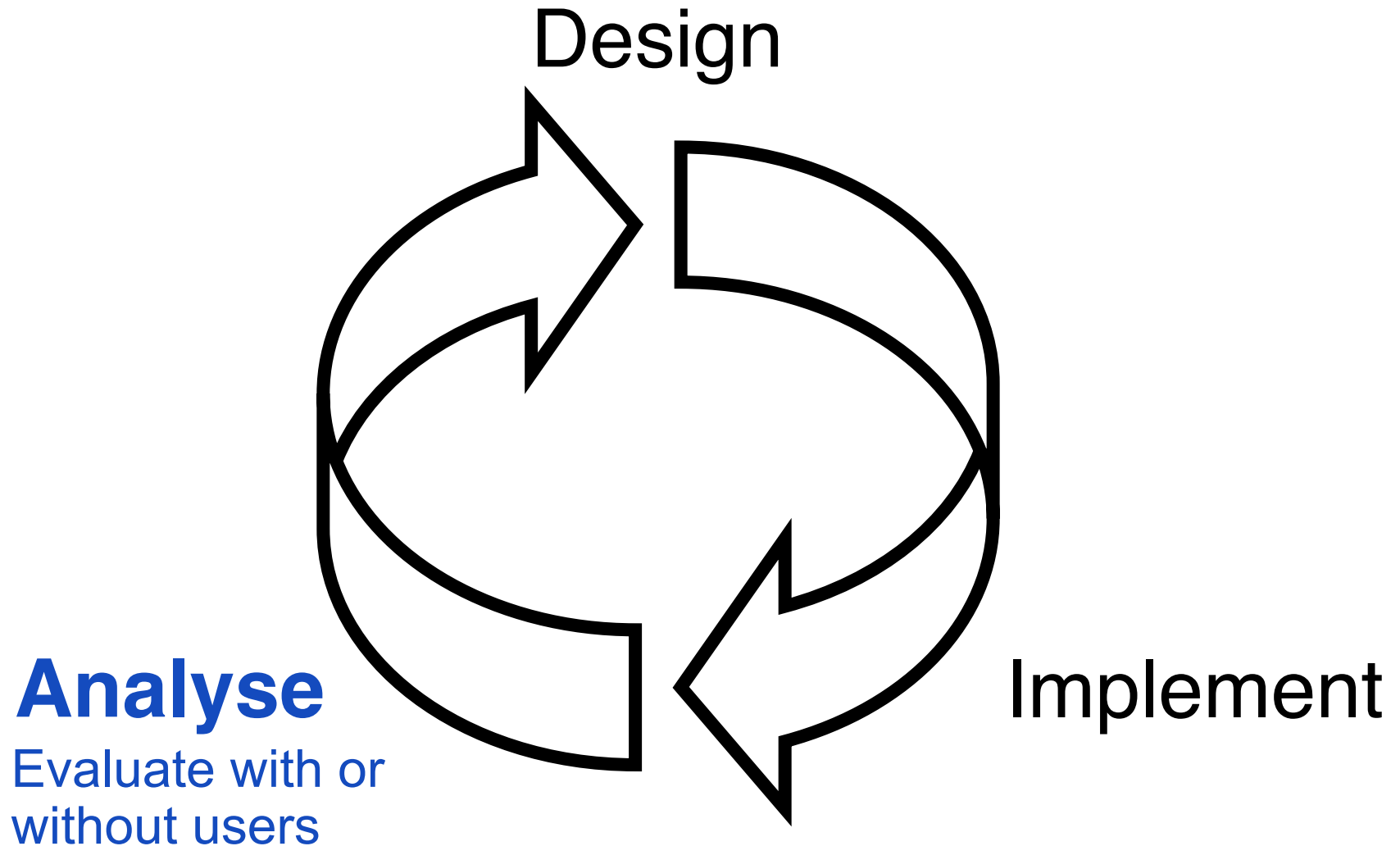
Chin, Diehl, Norman: Development of an instrument measuring user satisfaction of the human-computer interface. CHI '88



# EVALUATION



# DIA Cycle: When to evaluate?







## Where to evaluate: Laboratory



- + Equipment (audio / video, see-through mirrors, special computers), no disruptions, quiet
- Natural environment missing (shelves, wall calendar, streets, people...); unnatural situation (relevance?)

Only place possible if real use un-entical, remote (ISS...), or controlled situation needed



- Studies in the users' natural environment
- Advantages
  - + Situations (location and context!) and behavior more natural
  - + More realistic (also because of disruptions)
  - + Better suited to long-term studies
- Disadvantages
  - Noise, task interruptions, effort
  - Could still feel like a test situation





# Evaluation in the Mobile Context

- Context of use needs to be taken into account
  - Factors: User, activity, device, environment
- Usage “on the move”
  - Physically moving: walking, driving a car, traveling as a passenger
  - Being in different places: away from office environment or home
- Difficult to collect data in the field
  - Recording interaction
  - Capturing context
  - Controlling experimental conditions





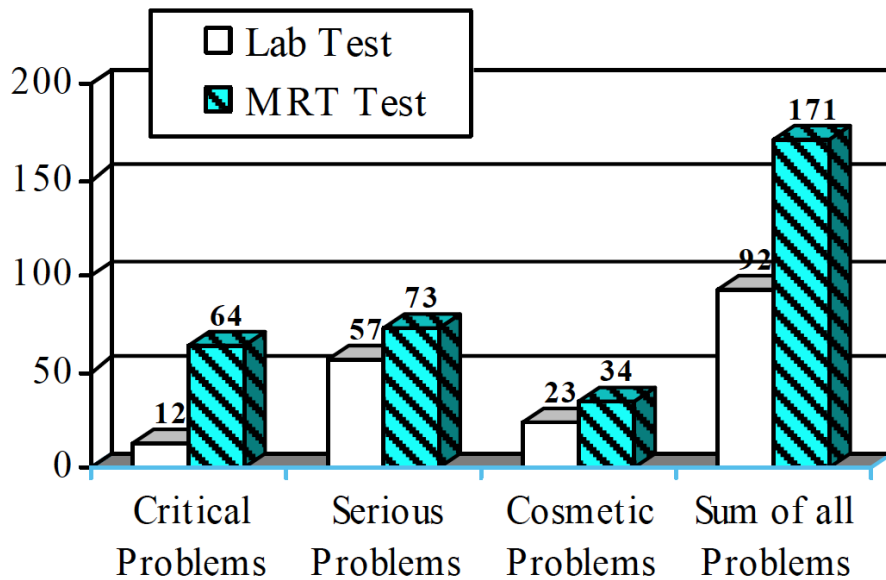
- Assess quantity and quality of usability problems found in lab vs. field
- Tasks and scenarios given

**Table 1. The Test Scenario and Tasks involved**

Task	Scenario of action	Task Description
1	You need to inform your friend about your personal particulars as he needs to fill up a form for you. You decide to call out.	<ol style="list-style-type: none"><li>1. Dial out to contact Gerald from mobile phone contacts list.</li><li>2. Start a conversation upon pick up as you normally would.</li><li>3. Verbally inform the contact your full name, NRIC, address and date of birth</li></ol>
2	You receive a call from a friend on your mobile phone. You answer the phone call.	<ol style="list-style-type: none"><li>1. Answer phone call as you normally would.</li><li>2. Start a conversation with the friend.</li></ol>
3	You need to inform your friend about your personal particulars information as he needs to fill up a form for you. You decided to SMS	<ol style="list-style-type: none"><li>1. Compose a SMS including the following information: your full name, NRIC, address and date of birth.</li><li>2. Send SMS to Gerald from mobile phone contact list.</li><li>3. Reply again to Gerald if necessary, i.e. if Gerald replied your message.</li></ol>

Image source: Duh, Tan, Chen: [Usability Evaluation for Mobile Device: A Comparison of Laboratory and Field Tests](#). MobileHCI 2006.

Problems found:



User behavior:

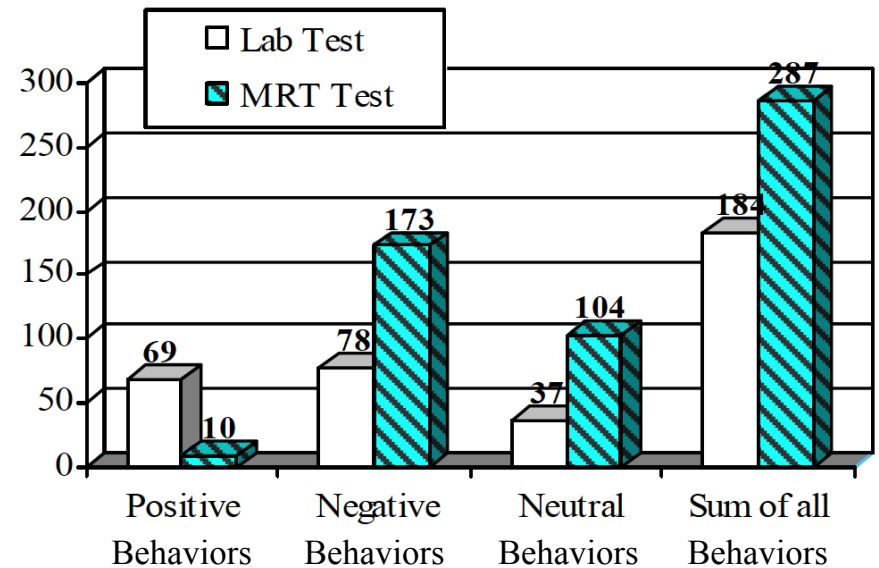
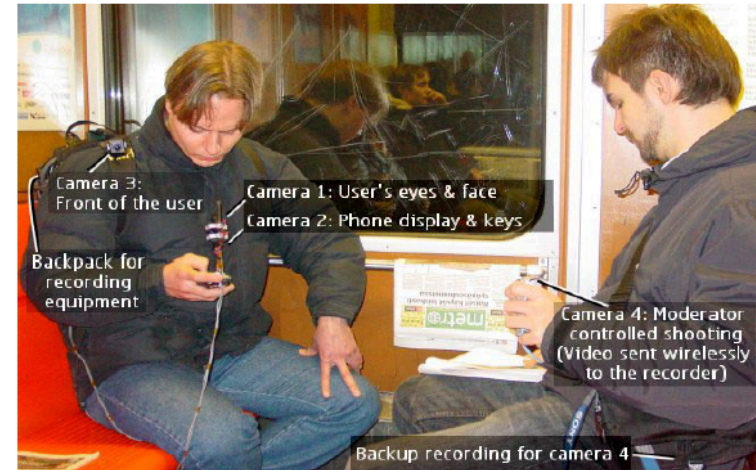


Image sources: Duh, Tan, Chen: [Usability Evaluation for Mobile Device: A Comparison of Laboratory and Field Tests](#). MobileHCI 2006.

- Evaluating the competition for cognitive resources when mobile
- Field study in urban environment
  - Performance of mobile Web tasks
  - Movement through urban situations
- Attention during loading a page
  - Duration of continuous attention
    - Lab: 16.2s → field: 4s
  - Number of attention switches
    - Lab: 1 → field: 8
  - Switching-back durations
    - Railway station: 7-8s, quiet street: 4-6s

Oulasvirta, Tamminen, Roto, Kuorelahti. [Interaction in 4-second bursts: the fragmented nature of attentional resources in mobile HCI](#). CHI '05.



**Figure 2. Configuration of recording equipment.**



**Figure 3. Output video data integrated on-the-fly.**



## Evaluating Without Users

- E1 Literature Review
- E2 Cognitive Walkthrough
- E3 Heuristic Evaluation
- E4 Model-Based Evaluation

## Evaluating With Users

### Qualitative

- E5 Conceptual Model Extraction
- E6 Silent Observation
- E7 Think Aloud
- E8 Constructive Interaction
- E9 Retrospective Testing

+ Interviews,  
questionnaires,...

### Quantitative

- E10 Controlled Experiments



# E1: Literature Review

- Many research results about user interface design have been published
- Idea: Search literature for evidence for (or against) aspects of your design
- + Saves effort, time and money by avoiding own experiments
- Results only carry over reliably if context (users, assumptions) is very similar





## E2: Cognitive Walkthrough

- Analytical method for early design or existing systems
  - Without users
- Expert evaluator = designer or cognitive psychologist
- Goal: Judge **learnability** and ease of use
  - Does system help user to get from goals to intentions and actions?
- Step through each action and ask
  - Is the effect of the action the same as the user's goal at that point?
  - Will users see that the action is available?
  - Once users find the action, will they know it is the right one?
  - After the action is taken, will users understand the feedback?



## E2: Cognitive Walkthrough

- What you need
  - Interface description (prototype of the system)
  - Task description
    - **Example:** Program the DVR to time-record a program starting at 18:00 and finishing at 19:15 on BBC 1 on June 2, 2018
  - List of interface actions to complete the task
  - User profile
- Doing the actual walkthrough
  - Analyze process of performing the actions using above questions
- Written questions capture psychological knowledge and guide the tester



# E3: Heuristic Evaluation

- Choose usability heuristics
  - (general usability principles, e.g., Nielsen's 10 Usability Principles)
- Step through tasks and check whether guidelines are followed
- Severity rating for each problem (Nielsen)
  - 0 = I don't agree this is a problem at all
  - 1 = cosmetic problem
  - 2 = minor usability problem, low priority to fix
  - 3 = major usability problem, high priority to fix
  - 4 = usability catastrophe, imperative to fix before release
- + Quick and cheap
- Subjective (have several independent evaluators)  
See also: [www.useit.com/papers/heuristic](http://www.useit.com/papers/heuristic)



# 10 Usability Principles (Jakob Nielsen)

1. Keep the interface simple!
2. Speak the user's language!
3. Minimize the user's memory load!
4. Be consistent and predictable!
5. Provide feedback!
6. Design clear exits and closed dialogs!
7. Offer shortcuts for experts!
8. Help to recover from errors, offer Undo!
9. Prevent errors!
10. Include help and documentation!



# 10 Rules for a good Crypto API? Smith & Green @ USENIX Hotsec'15

1. Easy to learn, **even without crypto background**
2. Easy to use, even without documentation
3. Hard to misuse. **Incorrect use should lead to visible errors**
4. **Hard to circumvent errors – except during testing/development**
5. Easy to read and maintain code that uses it
6. Sufficiently powerful to satisfy (**non-security**) requirements
7. ~~Easy to extend~~ **Hard to change/override core functionality**
8. Appropriate to audience – **this means people with no crypto experience**
9. **Assist with/handle end-user interaction**
10. **However, where possible integrate into standard APIs so normal developers never have to interact with crypto APIs in the first place**

conduct developer studies



- Principles, Heuristics
  - Small set of general rules (low authority, high generality)
  - Abstract rules, based on psychological knowledge
  - Largely independent of technology
- Guidelines
  - Large set of detailed rules (medium authority, low generality)
  - Often developed for a specific platform
  - More concrete, more technology-oriented
- Standards
  - Agreed upon by a large community (high authority, medium generality)
  - Carefully developed by a standards committee (consensus-based)



## Examples:

- ISO 9241: “Ergonomics of Human System Interaction”, 17 parts
  - 7 parts concerning hardware issues, 8 parts concerning software issues
- ISO 14915: “Software ergonomics for multimedia user interfaces”, 3 parts
  - “Multimedia navigation and control”, “Media selection and combination”



## Evaluating Without Users

- E1 Literature Review
- E2 Cognitive Walkthrough
- E3 Heuristic Evaluation
- E4 Model-Based Evaluation

## Evaluating With Users

### Qualitative

- E5 Conceptual Model Extraction
- E6 Silent Observation
- E7 Think Aloud
- E8 Constructive Interaction
- E9 Retrospective Testing

+ Interviews,  
questionnaires,...

### Quantitative

- E10 Controlled Experiments

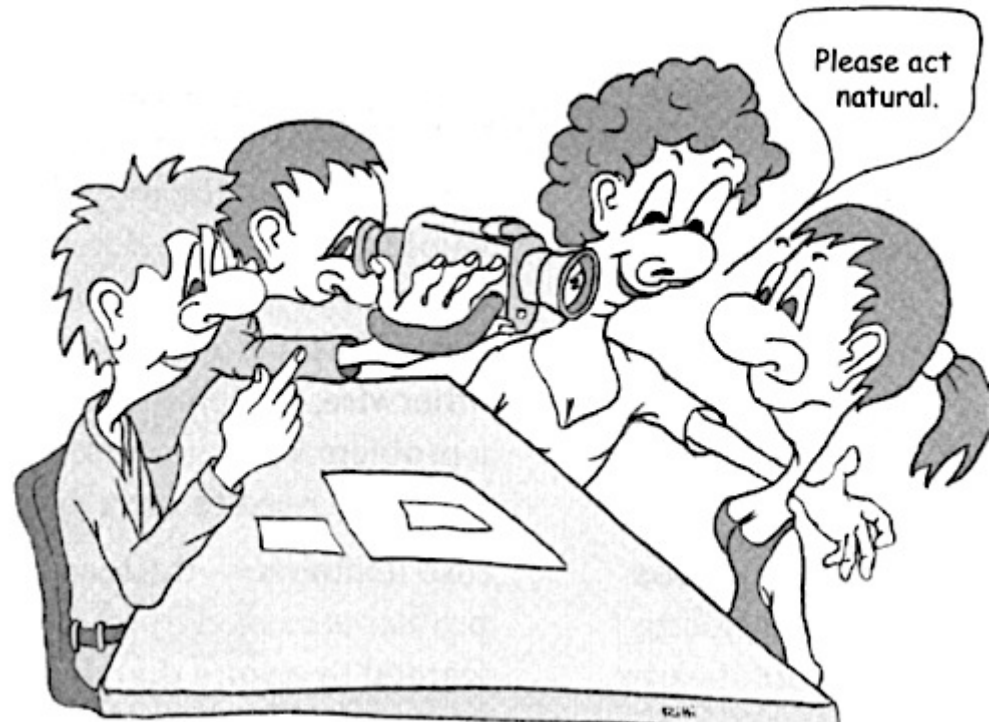




# Four Key Issues

1. Setting goals
  - Decide how to analyze data once collected
2. Relationship with participants
  - Clear and professional
  - Protect privacy
  - **Informed consent-form**
    - Signed agreement between evaluator and participant
3. Triangulation
  - Use more than one approach
  - Use different perspectives to understand a problem or situation
4. Iterate
  - If questions reveal that goal was not sufficiently refined: refine goal, repeat

- Tests are uncomfortable for the tester
  - Pressure to perform, mistakes, competitive thinking
- So treat testers with respect at all times!
  - Before, during, and after the test





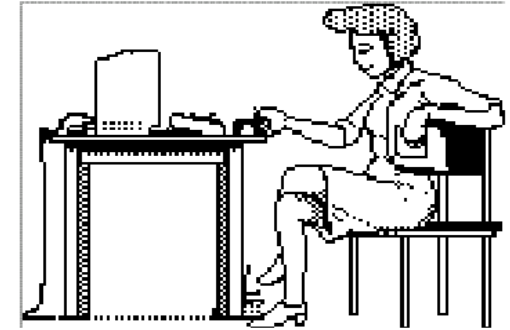
- Paper and pencil
  - Evaluator notes events, interpretations, other observations
  - Cheap but hard with many details (writing is slow)
  - Forms can help
- Audio recording
  - Good for speech with Think Aloud and Constructive Interaction
  - But hard to connect to interface state
- Video
  - Camera on user +
  - Screen-capture
  - Best capture, but may be too intrusive initially
  - Time consuming to process during evaluation
- Logging
  - Log input events of the user, synchronize with audio & video



- **Test**
- **Test**
- **And test again**



# Silent Observation



Source: Saul Greenberg

- Designer watches user in lab or in natural environment while working on one of the tasks
- No communication during observation
- + Helps discover big problems
- + Prevents “overly helpful” assistant problem
- No understanding of decision process (that may be wrong) or user’s mental model, opinions, or feelings
- Can end in almost zero result, if the participant gets things wrong in the beginning



# Think Aloud

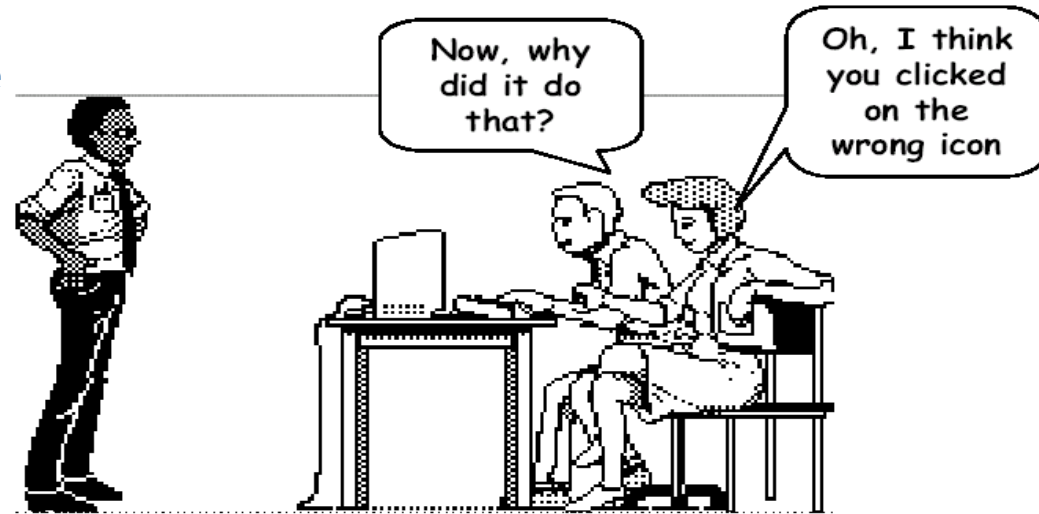


Source: Saul Greenberg

- As Silent Observation, but user is asked to say aloud
  - What he thinks is happening (state)
  - What he is trying to achieve (goals)
  - Why he is doing something specific (actions)
- Most common method in industry
- + Good to get some insight into user's thinking, but:
  - Talking is hard while focusing on a task
  - Feels weird for most users to talk aloud
  - Conscious talking can change behavior



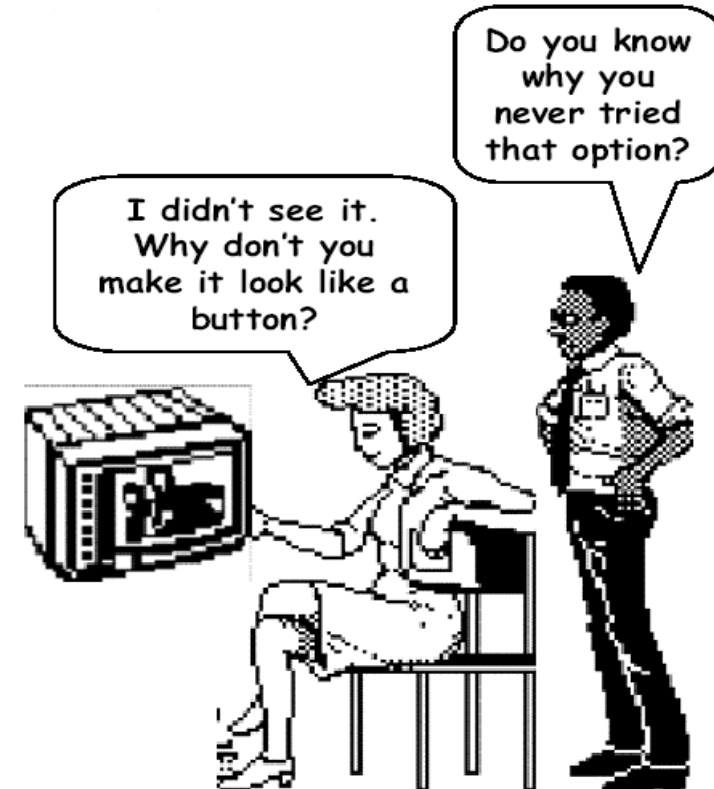
## E8: Constructive Interaction



- Two people work on a task together
  - Normal conversation is observed (and recorded)
    - + More comfortable than Think Aloud
- Variant of this: Different partners
  - Semi-expert as “trainer”, newbie as “student”
  - Student uses UI and asks, trainer answers
    - + Gives insight into mental models of beginner and advanced users at the same time!



# E9: Retrospective Testing



- Additional activity after an observation
- Subject and evaluator look at video recordings together, user comments his actions retrospectively
- Good starting point for subsequent interview, looking at video avoids wrong memories
- Often results in concrete suggestions for improvement





# INTERVIEWS



- + Ability to go deep
  - by asking that explore a wide range of concern
  - and giving interviewees the freedom to provide detailed responses
  - data can be gathered that would otherwise be very hard to capture
- This flexibility comes at a price:
  - potentially unbounded discussion must be managed
  - interviews are harder to conduct than surveys
  - time consuming
- Analysis is a major challenge
  - transforming and merging raw notes into usable data is challenging and time consuming
- Interview is separate from task
  - problem of recall
  - self-reported data



# The Interview

- Tell the interviewee that he/she can decline to answer any question or opt-out entirely at any time
- State what the study is about (+/- subterfuge)
- Start with relatively easy questions
  - this builds trust and confidence
- Intersperse hard questions with easy questions
  - to defuse any tension
- Critical questions should be done near the end
- End the interview with easy questions
  - this creates the feeling of accomplishment
- Debrief the interviewee and/or offer a follow-up information exchange
- Thank the interviewee for their time and participation



# Debriefing

- Turn of recording devices for debriefing
  - participants might share comments they were not willing to say during the interview
  - care is needed in dealing with this data!



# Try it out

- Mini interview
- Please form groups of two
  - one interviewer
  - one interviewee
- Task:
  - conduct an interview to learn about the technical procedure how to construct a peanut butter and honey sandwich
  - assuming that you have new jars of peanut butter and honey, a fresh loaf of bread and a standard kitchen
  - **find as many technical details as possible!**
  - use your smart phones to record the interview (optional)



# Discussion

- Get plate
- Get cutlery
- Open bread packaging
  - were you asked if the bread came in packing?
  - did you have to decide on the fly?
- Cut bread with bread knife
- Open jars
- Remove foil
  - this is often forgotten since it is implicit knowledge
- Spread peanut butter and honey
  - preferences for order?
  - spoon for the honey?
  - same knife for both?
  - single slice or double decker?



# Types of Interview

- Unstructured / Open ended / Exploratory
  - Some specific questions that are planned
  - But based on interviewees responses interviewers can
    - re-order questions
    - invent new lines of inquiry
- Structured
  - rigid script with the questions
  - questions can be skipped based on pre-defined rules
  - similar to survey
    - but it might be easier to answer a question verbally than to write it down
  - easier to analyse
- Semi-Structured
  - mix of the above



- Closed questions
  - yes-no
  - multiple choice
  - likert scale
  - Do you like the design of this warning message?
- Open questions
  - Questions asking for responses, opinions, feedback without external constraints
  - What did you think about the design of this warning message?
- Tasks (hidden questions)
  - Please complete this sentence:
    - “The most frustrating problem with PGP is...”
- Conceptual Mapping
  - “Please draw a diagram showing how web security works”





- Avoid compound questions
  - What were the strengths and weakness of the menu and the toolbar?
  - + What did you think about the menu layout?
  - + What did you think of the toolbar?
  - + Which did you prefer?
  
- Paraphrase complex or unclear answers
  - “Did I understand you correctly, you think that...”
  - The closed yes/no questions allow you to be sure that you extracted the right information



# Try it out

- Mini interview
- Please form groups of three (7 min., then change roles)
  - one interviewer
  - one interviewee
  - one referee
- Task:
  - Find out the details of a problem of choice of the interviewee
  - Examples
    - Problems of using PGP
    - Trust in Crypto-currencies
  - Only use open questions!
  - If closed question is asked interviewee answers only with yes or no and says nothing else
  - Referee buzzes when these rules are broken
  - Use paraphrases to check if you understood correctly



- Questions should be as unjudgmental and unbiased as possible
- Watch out for phrasing that could encourage participants to give answer they think you want to hear
  - Particularly if you are asking about something you built
- Examples:
  - Why do you like this design?
  - Don't you think this is difficult to use?
  - Did you like...?
  - + What did you think of...?



## Evaluating Without Users

- E1 Literature Review
- E2 Cognitive Walkthrough
- E3 Heuristic Evaluation
- E4 Model-Based Evaluation

## Evaluating With Users

### Qualitative

- E5 Conceptual Model Extraction
- E6 Silent Observation
- E7 Think Aloud
- E8 Constructive Interaction
- E9 Retrospective Testing

+ Interviews,  
questionnaires,...

### Quantitative

- E10 Controlled Experiments



# E10: Controlled Experiments

- Quantitative, empirical method
- Steps
  - Formulate hypothesis
  - Design experiment, pick variable(s) and fixed parameters
  - Choose subjects
  - Run experiment
  - Interpret results to accept or reject hypothesis



- Statistical Unit
  - Objects from which measurements are collected (e.g. a participant or a country)
- Population
  - The set of all statistical units relevant to a particular investigation
- Sample
  - The subset of the population that was actually analysed
- Attribute/Parameter/Variable
  - A property of the statistical unit that we are interested in
- Value
  - The actual data for a variable measured for one statistical unit



# Two Main Types of Variables

- Independent (IV) / Predictor / Factor / Input
  - what we base our explanation on
  - characterize statistical units
  - examples: age, expertise
- Dependent (DV) / Outcome / Target / Output
  - what we are trying to explain
  - examples: usability, time taken



- A claim that predicts outcome of a DV based on an IV
- Approach: Reject null hypothesis (inverse, i.e., “no influence”)
  - Null hypothesis is a term from statistical testing
- Consider this hypothesis:

*“There is no difference between the target selection speed when using a mouse, a joystick, or a trackball to select icons of different size (small, medium, large)”.*
- What are the DV and the IVs? How many conditions in the experiment?

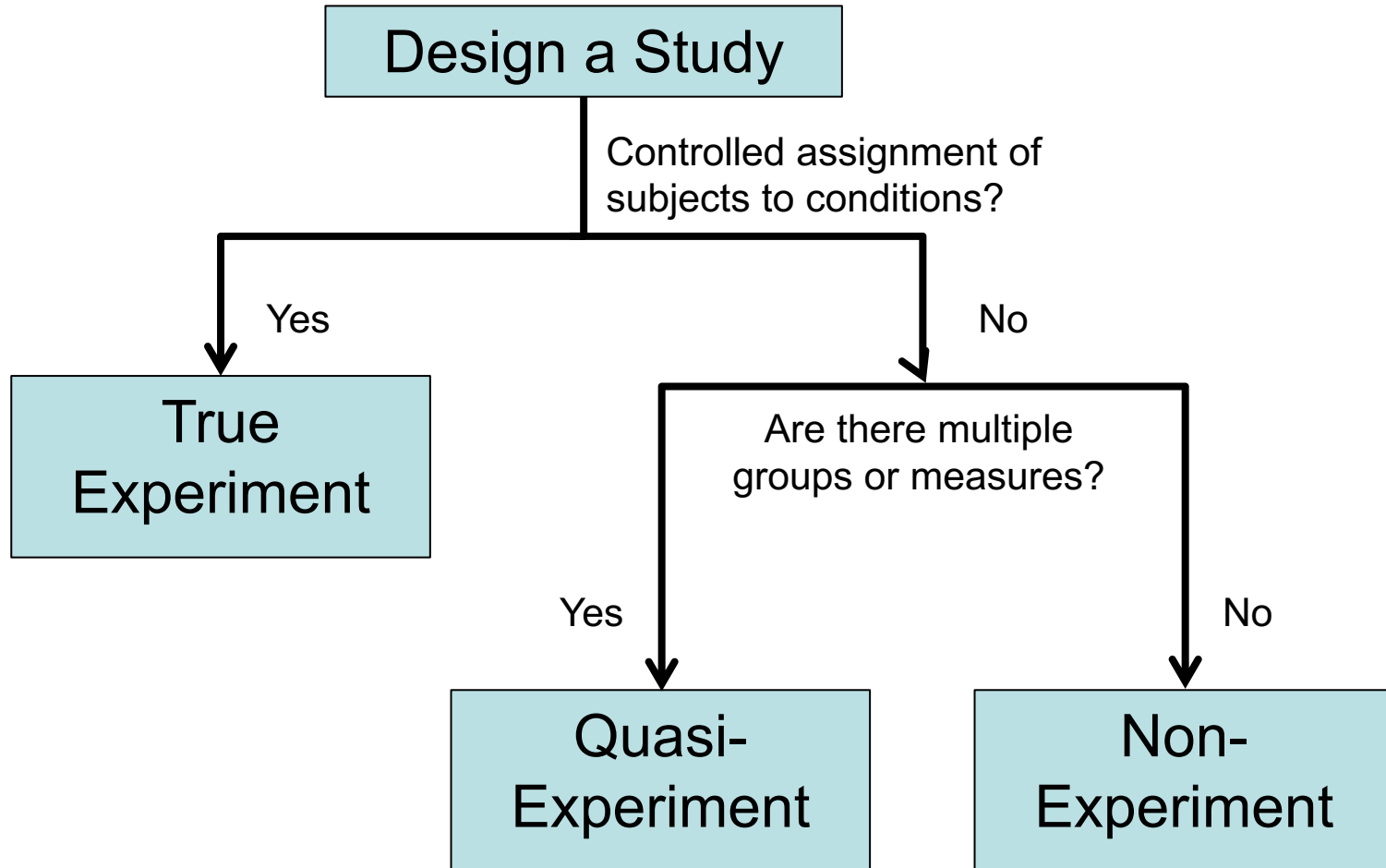




- Condition
  - The procedure that is varied in order to estimate a variable's effect by comparison with a control condition
- Number of conditions
  - product of the number of values in each IV
  - in our example:  $3 \times 3 = 9$
- How to control the IVs or condition assignment?
  - straightforward in previous example: get joystick, mouse or trackball
  - challenging in other cases:
    - testing influence of demographic properties (cognitive abilities, left-handed vs. right-handed, ...)
    - testing against best-case scenarios that do not exist (Wizard of Oz).



# Types of Research Designs





- True Experiment
  - the experimenter controls assignment of experimental units (e.g., participants, rats) to experimental conditions
  - control allows to draw causal conclusions
  - example:
    - randomly give half of the participants the real drug and the other half a placebo
    - let half of the participants randomly use one product or the other and measure usability.

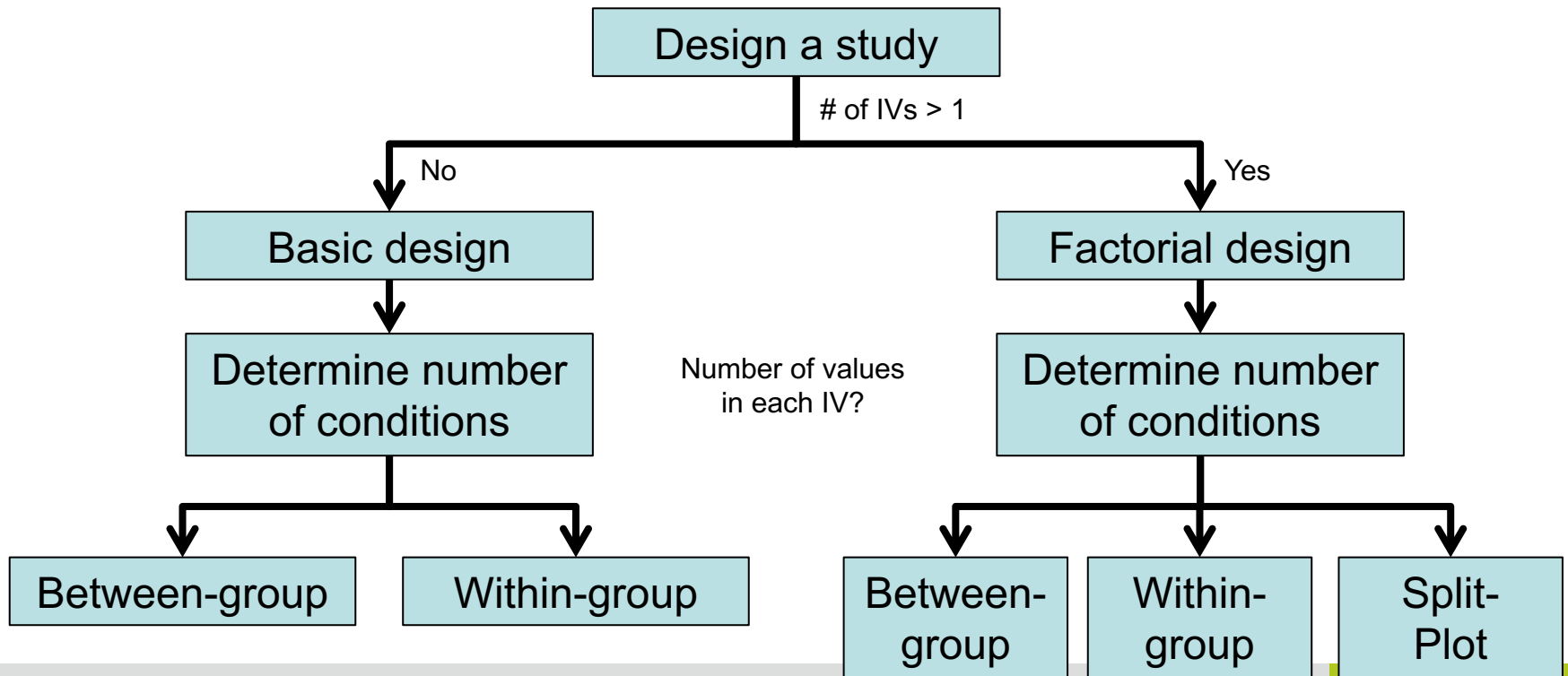


- Quasi-Experimental
  - if assignment cannot be controlled, another assignment criterion is used
  - examples:
    - some statistical units occurred before a certain “*external*” event and some after (movie revenue before and after Megaupload was closed down)
    - compare pupils according to grade averages above and below a threshold or from one type of school to another
    - compare people who have a security background with those who don't
  - Remember: more than one condition is needed to be a quasi-experiment



- Non-Experimental
  - describe phenomenon “as is”
  - do not manipulate variables
    - therefore, cannot deduct a cause!
  - e.g. surveys, ethnography(, interviews), ...
  
- These three forms can be mixed, especially when conducted online.

- Goal: draw a big picture of how to run the experiment
  - estimate a timeline (and a budget)
- Two essential questions:
  - How many IVs do we want to investigate in this experiment?
  - How many different values does each IV have?





- Consider the following hypotheses. How many conditions in each hypothesis?
  - H1: “There is no difference in typing speed when using a QWERTY keyboard, a DVORAK keyboard, or an alphabetically ordered keyboard.”
  - H2: “There is no difference in the time required to locate an item in an online store between novice users and experienced users.”
  - H3: “There is no difference in the perceived trust towards an online agent among novice and experienced customers who are from the United States, Russia, China, and Nigeria.”



- How many conditions do we expose each participant to?
  - Between-group / between-subject design
    - The effect of each condition is measured between groups/subjects.
      - i.e. each participant is exposed to one condition only.
    - If the task is to type 500 words using a selected keyboard, each participant types 500 words.
  - Within-group / within-subjects design
    - The effect of each condition is measured within the group/each subject.
      - i.e. one group of participants is exposed to all conditions
    - In this case, each participants types 1500 words.
- This decision implies the use of different statistical analyses.





# Active Learning

- Think – Pair – Share
- What are the advantages/disadvantages/differences between
  - Within subjects
  - Between subjects
- ?



- Cleaner design
  - no learning from previous exposures
  - less time spent in the experiment
    - ➔ less influence of fatigue and frustration
  
- Compare two distinct groups of participants
  - there is no baseline for every individual
    - individual differences cause noise
    - need to make sure groups are very similar
  - number of participants in each group needs to be high
    - sample size = no. of conditions x no. of participants per condition (as dictated e.g. by power analysis).
    - example: 4 conditions x 16 participants/condition:  $N=64$



# Within-group Design

- Exactly the opposite of between-groups design:
  - smaller sample size needed:  $N=16$  (!)
    - therefore usually a lot cheaper to conduct
  - more meaningful measurements due to individual comparisons
  
- Disadvantages
  - repeated exposure can cause learning and fatigue
    - learning favours subsequent exposures
    - fatigue favours initial exposures
    - can add to substantial overall bias



- Difficult decision, to be made on case-by-case basis
- Hybrid setups possible (later slides)
- Between-groups should be used for:
  - simple tasks with limited individual differences
    - limited cognitive processes
    - e.g. basic motor skills when selecting a screen target
      - as opposed to reading, comprehension, information retrieval and problem solving
  - tasks that would be greatly influenced by learning effects
    - first-contact required, e.g. when testing website design
  - problems that cannot be investigated using within-groups design
    - consider H2 and H3 from before



# Using Between-group Design

- Randomly assign participants to conditions
  - randomly does not mean haphazardly!
  
- Counterbalance confounding factors
  - gender
  - age
  - computing experience
  - Internet experience
  - ...
  - i.e. all (“relevant”) demographic properties except those that are IVs
  
- Make sure groups are as similar as possible w.r.t. your hypothesis.



# When to use Within-group Design

- Within-groups design isolates individual differences more effectively
- Within-groups should be used for:
  - tasks with large individual differences
    - i.e. reading, comprehension, information retrieval and problem solving
  - tasks that are less susceptible to learning effects
    - involving complicated or learnt skills and knowledge
    - for example: investigating the impact of font type on reading speed
  - very small target participant pools
    - for example, when looking for a particular participant property (disabilities, illnesses, or a combination of demographical properties)



- Need to control for negative impact of learning, fatigue and other within-groups problems
  - randomise task order to control for learning
    - for example, learning effects in one participant using the DVORAK keyboard last are offset by another participant using it first
  - provide a training session
    - if participants can familiarise with all conditions before the actual experiment, learning has less influence
    - commonly used in combination with task randomisation
  - limit the total time spent in the experiment
    - generally between 60 and 90 minutes or less
    - never more than 120 minutes
    - plan (force) breaks when necessary

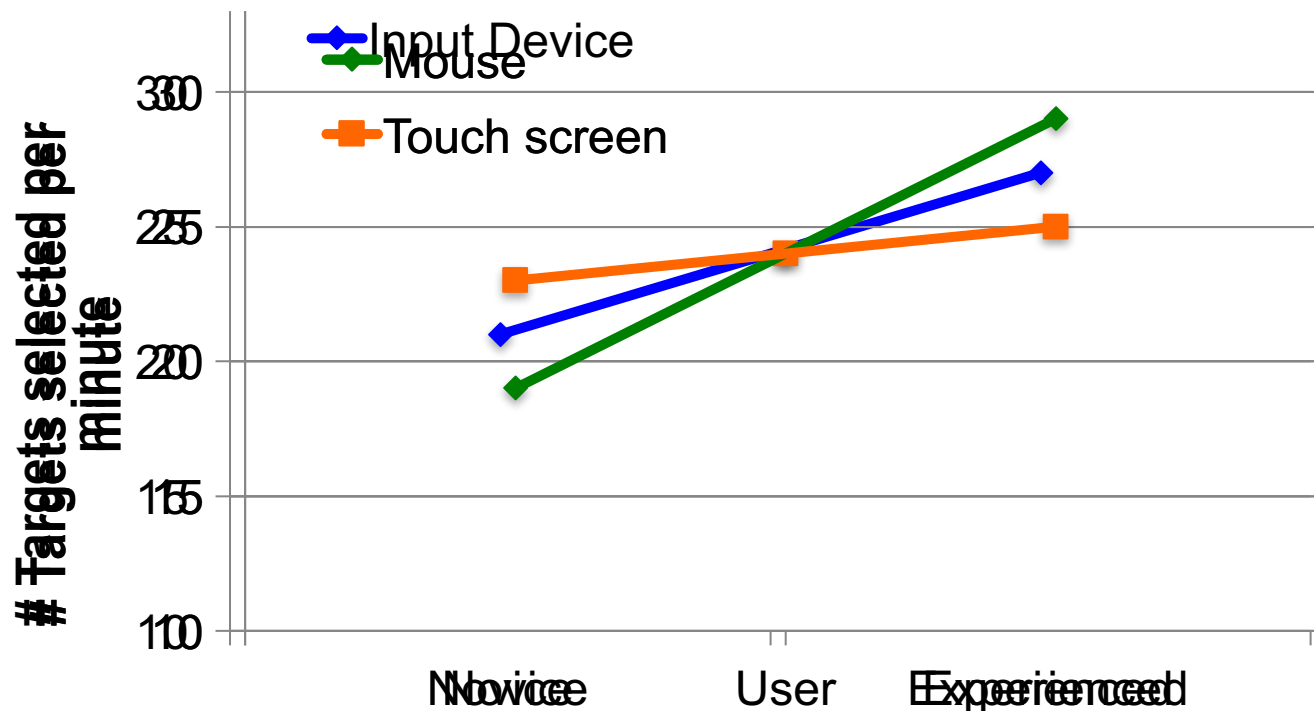


- Used to investigate more than one IV
  - number of conditions is the product of the number of values in each IV
  - example: in addition to three **keyboard types** (QUERTY, DVORAK, and Alphabetic), we also want to investigate the effects of **different tasks** (composition vs. transcription) on typing speed.
  - $3 \times 2 = 6$  conditions
    - across the **task dimension**, we can examine the impact of keyboards
    - across the **keyboard dimension**, we can examine the impact of task type
    - across **both IV dimensions**, we can examine interaction effects





- “The effect of one IV on the DV, depending on the particular value of another IV”.
  - one IV alone may not cause significant effects
  - interaction effects can provide additional insights
- Example (IVs: input device and expertise, DV: selection speed)

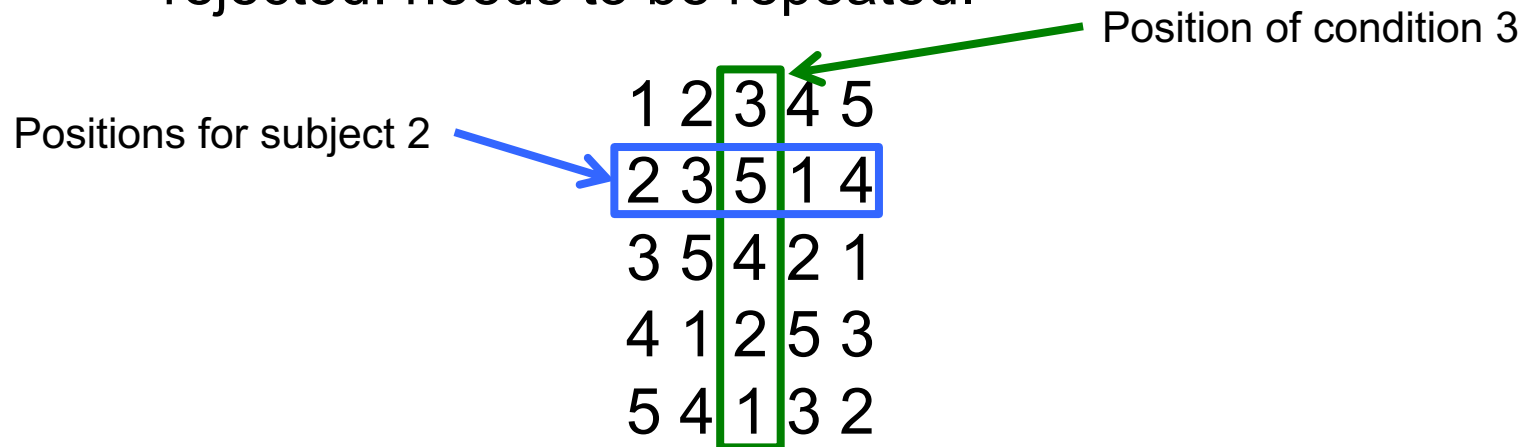




- Investigate some variables within-groups and others between-groups
  - example: testing the influence of age and GPS assistance to driving efficiency.
    - age cannot be tested within-groups
    - but each driver can drive with and without GPS assistance
  - advantage: smaller sample but still baseline for each participant in some IVs



- Variation of Split-Plot: Latin Squares
  - Latin squares make sure that each condition is assigned to each position the same number of times in a within-subjects study.
  - simulates a between-subjects design on each position
  - is not a truly random assignment!
  - can be difficult to administer
    - Latin square property is broken when one record is rejected: needs to be repeated.





- Reliable Experiments can be replicated by others teams, in other locations and another time
  - yielding results that are consistent, dependable, and stable
- Big challenge for HCI compared to “hard sciences”
  - measurements of human behaviour and social interaction subject to higher fluctuation and individual differences
  - less replicable
  - especially when using smaller samples
- Fluctuations in experimental results: Errors
  - okay: random errors
  - bad: systematic errors!



- If observing someone typing documents in five sessions, we may obtain a typing speed of 50 words per minute
  - in each session, speeds will vary:
    - 46 words/min, 52 w/m, 47 w/m, 51 w/m, 53 w/m
    - on average: about 50 words/min
- Observed values = actual value + random error
  - random error = noise
  - occurring by chance
  - push actual value up and down
  - the larger the sample, the more the noise equalises
  - therefore, can be ignored



# Systematic Errors

- Non-random error
  - also called bias
  - systematically influencing measurement values
    - pushing values into the same direction
  - no offsetting in larger samples
  - typing example:
    - imagine the experimenter is looking over the secretary's shoulder and then gets these results:
      - 47,44,45,42,46 → mean=44.8
    - still fluctuating, but underperforming



# Systematic Errors

<u>Instability</u>	<u>Biasedness</u>		
	Much	Some	None
Much			
Some			
None			

- Instability of sampling vs. bias (Rosnow & Rosenthal, 2008)
- Bias is the true enemy of experimental research
  - try to eliminate and control bias
  - isolate impact when inevitable
  - five major sources



- Instruments used to measure values can introduce systematic error
  - e.g. ill-calibrated sensor, slow stopwatch, misused questionnaire
- Can be avoided by carefully selecting and examining instruments before the experiment





- Inappropriate or unclear experimental procedures may introduce biases
  - learning, fatigue
  - wording of instructions
    - need to be consistent across conditions and instructors!
    - example: participants instructed to work “*as quickly as possible*” were slower than those instructed to “*take your time, there is no rush*”.
  - ambiguous instructions can cause unwanted variations:
    - e.g., some participants may hold a smartphone sideways and be quicker than those who hold it regularly
    - can be avoided by clarifying instructions
- Pre-testing is essential!



- To reduce bias through experimental procedures
  - randomise the order of conditions, tasks, and task scenarios in within-group and split-plot designs
  - prepare a written document with detailed, tested and revised instructions for participants
  - run pilot studies/pre-tests with actual participants before collecting real data and revise your procedure accordingly
  - prepare a written document with detailed, tested and revised instructions for experimenters
    - even if you are the only experimenter!



# Participant and Recruiting Bias

- Consider participant demographics and possible biases
  - e.g., recruiting IT students only will bias your results!
- Recruit carefully: make sure your participant pool mirrors the target user population (as best as possible)
- Create an environment and task procedure that causes the least stress to the users
- **Reassure participants that you are testing a product/an interface/a method and not them.**
  - has been shown to make people more calm and relaxed during experiments
- Be prepared to reschedule exhausted, stressed, tired or very nervous participants.



- Major source of bias
  - may intentionally or unintentionally influence experiment results
    - worse if multiple experimenters are used: different influences are hard to discover and quantify
  - spoken language, body language, and facial expressions cause bias!
    - “I think you will like this... I designed it myself!”
    - a frustrated looking experimenter says: “Damn! That’s so slow!”
    - a waiting experimenter leans forward and taps his fingers on the mouse while waiting for an application to load.
    - the experimenter arrives late and after the participant. He then takes 10 minutes to set everything up...



- Controlling for experimenter behaviour bias:
  - train experimenters or practice yourself
    - stay neutral, calm and patient during experiments no matter what happens
  - experimenters should make sure to arrive 15 minutes ahead of scheduled sessions to prepare everything.
  - supervise a session using two experimenters when possible.
    - lead experimenter interacts with participants, assistant observes and takes care of technical problems.
  - prepare written documents with detailed procedures that all experimenters have to follow strictly.
  - when appropriate, record instructions for the participants before the experiment and play the recording during the experiment.



- Make sure the environment is agreeable
  - quiet room, appropriate lighting, comfortable chairs and tables, clean, tidy, no distractions.
  - if possible, the participant should be seated alone and the experimenter observes from another room.
  - in a field study, check the study location before the scheduled time.
- A final note on bias: we can try to avoid bias, but there always will be some bias.
  - We should be careful when interpreting and reporting the findings.
  - We should be ready to argue why a certain bias is acceptable.
    - for example: using only students



- Between-groups
  - Each subject only does one variant of the experiment
  - There are at least 2 variants (manipulated form & control, to isolate effect of manipulation)
  - + No learning effect across variants
  - But requires more users
- Within-groups
  - Each subject does all variants of the experiment
  - + Less users required, individual differences canceled out
  - But often learning effect across variants problem



- Statistical analysis
  - Often assumptions about underlying distribution
  - t-test: Compare two groups, normal distribution
  - Analysis of variance (ANOVA): Compare two or more groups, normal distribution
  - Regression analysis: How well does result fit to a model?
  - Wilcoxon- or Mann/Whitney test,  $X^2$  test
- Choice depends on
  - Number, continuity, and assumed distribution of dependent variables
  - Desired form of the result (yes/no, size of difference, confidence of estimate)





# Other Evaluation Methods

- Before and during the design, with users
  - Personal interviews
  - Questionnaires
- After completing a project
  - Email bug report forms
  - Hotlines
  - Retrospective interviews and questionnaires
  - Field observations (observe running system in real use)



## Evaluating Without Users

- E1 Literature Review
- E2 Cognitive Walkthrough
- E3 Heuristic Evaluation
- E4 Model-Based Evaluation

## Evaluating With Users

### Qualitative

- E5 Conceptual Model Extraction
- E6 Silent Observation
- E7 Think Aloud
- E8 Constructive Interaction
- E9 Retrospective Testing

+ Interviews,  
questionnaires,...

### Quantitative

- E10 Controlled Experiments



# Summary

- Evaluate to ensure system matches users' needs
- Evaluation should happen throughout the design process
  - By experts (analytically)
  - By users (experimentally)
- A plethora of methods to evaluate designs
  - Decide when to apply which
- Treat testers with respect at all times!