# Usability Testing

Refs: Lecture note from D. Petkovic

# Objective

Ÿ In order to help you with designing and executing usability testing, I am summarizing main points from class 3 and adding a few more practical guidelines.

Ÿ Recommend book:

› Technical Communication Library, 1994

# Usability Testing

Ÿ Process that employs representative target population participants to evaluate product usability using specific usability criteria

Ÿ Usability testing is not a guarantee for product success, but it should identify at least the key problems

# Basic components

1. Development of specific problems statements and test plans and objectives
2. Use of representative sample of end users
3. Representation of the actual work environment
4. Observation of end users during product use or review
5. Collection of quantitative and qualitative measurements
6. Analysis of results and recommendations

# Types of usability tests

## *Exploratory:*

Ÿ Early in the process

Ÿ Can be based on any form of the GUI (sketch, wire diagrams etc.)

Ÿ Evaluate preliminary, basic design concept

Ÿ 

Ÿ Informal test methodology, a lot of interaction

Ÿ Discuss high level concepts

# Types of usability tests

*Assessment*

Ÿ   Done After fundamental concepts are done

Ÿ   Evaluates usability of lower level operations

Ÿ   The users actually perform set of well defined tasks

Ÿ   Less interaction with test monitor

Ÿ   Quantitative measurements are collected

# Types of usability tests
## *Validation*

Ÿ    Done late in development cycle, close to release

Ÿ

Ÿ    Often the first time when the whole product si tested (including help and docs)

Ÿ    Evaluate wrt. some predetermined usability standard or benchmark

Ÿ    Standards come from previous testing, competitive information, marketing etc.

Ÿ    Very specific quantitative tests

Ÿ    Can establish standards for future products

Ÿ    Can be also done by beta customers

# Types of usability tests

## *Comparison*

Ÿ Can be done at any point in the development cycle

Ÿ Compare alternatives using objective measures

Ÿ Can be informal of formal, depending when it is done

Ÿ Often, best of alternative designs is combined

# Test environments

- Ÿ Simple single room setups
  - › Observer/monitor close to evacuator
  - › Observer removed from evaluator
- Ÿ Electronic observation room
- Ÿ Classic elaborate usability lab
- Ÿ Mobile lab

# Typical test plan format

Ÿ  <u>Purpose</u>: what is the main purpose of the test

Ÿ  <u>Problem statement</u>: specific questions you want resolved

Ÿ  <u>Test plan and objectives:</u> tasks the user will do

Ÿ  <u>User profile:</u> who will be the users

Ÿ  <u>Method and test design:</u> how will you observe it, how will you collect the data etc.

Ÿ  <u>Test environment and equipment</u>

Ÿ  <u>Test monitor role</u>

Ÿ  <u>Evaluation measures and data to be collected:</u> how will you collect the feedback and how will you evaluate it

Ÿ  <u>Report:</u> what will final report contain

# Task selection for evaluation

Ÿ Tasks to be evaluated are functions users want to do with the product. Focus is on user view of the tasks and NOT at the components and details that you used to implement it. Examples:

› Create and file document

› Import several images

› Find the right document

Ÿ Objective is to indirectly expose usability flaws by asking the user to perform typical tasks and NOT telling them exactly how to do it.

Ÿ Choose key and most frequently done tasks

Ÿ The task has to be specific and measurable (quantitatively or qualitatively)

# Task components

| TASK | DESCRIPTION |
|---|---|
| Task | Load paper into copier |
| Machine state | Paper tray empty |
| Successful completion criteria | Paper properly loaded |
| Benchmark | Completed in 1 minute |

# Selection of evaluators and test groups

Ÿ Evaluators should be representative of the targeted users

Ÿ Independent groups or within-subject design (but be careful to avoid exposing users to same tests since this will bias the results)

Ÿ Adequate numbers of testers

Ÿ Offer motivation and rewards

# Measurements and Questionnaires

Ÿ <u>Performance data:</u> measures of user behavior such as error rates, number of accesses to help, time to perform the task etc.

  › Usually can and should be objectively and automatically measured

Ÿ <u>Preference data:</u> measures of user opinion, thought process such as rankings, answers to questions, comments etc.

  › Use questionnaires.

# Some performance measures (measure what can be measured)

- Ÿ Time to complete each task
- Ÿ Number and percentage of tasks completed successfully/unsuccessfully
- Ÿ Time required to access information
- Ÿ Count of incorrect selections
- Ÿ Count errors
- Ÿ Time for system to respond
- Ÿ

Data should be collected automatically or manually in an objective way.

# Questionnaires (for preference data)

*Likert scale*

Ÿ I found GUI easy to use (check one)

__ Strongly disagree  __ Disagree

__ Neither agree or disagree

__ Agree    __ Strongly agree

(can also assign numbers from  2 to 2)

*Semantic differentials*

Ÿ I found File Open menu (circle one)

Simple   3 2 1 0 1 2 3   Complex

# Questionnaires

## *Fill in questions*

Ÿ I found the following aspects of GUI particularly easy to use (list 0-4 aspects)

   --------------------------

   --------------------------

   --------------------------

   --------------------------

# Questionnaires

*Check-box*

Ÿ   Please check the statement that best describes your usage of spell check

      --- I always use spell check

      --- I use spell check only when I have to

      --- I never use spell check

# Questionnaires

## *Branching questions*

Ÿ Would you rather use advanced search

      --- NO (skip to question 19)

      --- YES (continue)

Ÿ What kind of advanced search would you like? (check one)

      --- Boolean

      --- Relevance

# Summarizing performance results

Ÿ Performance data

  › Mean time to complete

  › Median time to complete

  › Range (high and low)

  › Standard deviation of completion times

  › System response time statistics

Ÿ Task accuracy

  › % of users completing the task within specified time

  › % of users completing the task regardless of time

  › Same as above, with assistance

  › Average error rate

# Summarizing preference results

Ÿ **For limited choice questions**

› Count how many participants selected each choice (number and %)

› For Likert scale or semantic differentials provide average scores if there are enough evaluators

Ÿ **For free form questions**

› List questions and group answers into categories, also into positive and negative answers

Ÿ **For free comments**

› List and group them at the end of the report

# Analyzing Data

Ÿ   Identify and focus on tasks that did not pass the test or showed significant problems

Ÿ   Identify user errors and difficulties

Ÿ   Identify sources of errors

Ÿ   Prioritize problems by criticality = severity AND probability of occurrence

Ÿ   Analyze differences between groups (if applicable)

Ÿ   Provide recommendations at the end

# Problems statements and performance data to collect

| Problem Statement | Performance Data Collected |
|---|---|
| How effective is the tutorial | Compare error rates of users who used and not used this |
| How easy is it to perform task X | Error rate OR Number of steps needed |

Note: this is Performance data measurement only. You also need to asses user Preference data (see next slide)

# Problems statements and preference data to collect

| Problem Statement | Preference Data Collected |
|---|---|
| How effective is the tutorial | Ask user to rate it from very ineffective to very effective (Lickert scale or semantic differentials) + free comments |
| How easy is it to perform task X | Ask user to rate it from very easy to very difficult (Lickert scale or semantic differentials) + free comments |

# Relate problem statements with tasks

| Problem Statement | Task |
|---|---|
| How effective is the tutorial | GroupA: Import image w/o using tutorial<br><br>GroupB: Same but use tutorial first |
| How easy is it to Create Virtual Machine | Create Virtual machine with "this" properties using New VM Wizard |

# Task components

| TASK | DESCRIPTION |
|------|-------------|
| Task | Create VM using New VM Wizard |
| Machine state | VMware WS SW just loaded |
| Successful completion criteria | Working VM created |
| Benchmark | Completed in 30 sec. |

# Example questionnaire

Question 15:

It was easy to create virtual machine

◉          ○          ○          ○          ○

Strongly          Disagree          Neutral          Agree          Strongly
disagree                                                                agree

Comments:

# Some Suggested GUI issues to cover in preference data collection

Use GUI principles as general measures of quality to evaluate

Ÿ      Screen layout matches tasks

Ÿ      Amount of information is adequate

Ÿ

Ÿ      Good use of colors

Ÿ      Proper grouping of related info

Ÿ      Navigational problems

Ÿ      Users get lost in the system

Ÿ      Organized by user tasks

Ÿ      Icons are self-explanatory

Ÿ      Consistency

Note: ask these questions in the context of concrete user tasks not in vacuum