



Aalto University  
School of Arts, Design  
and Architecture

# Design B + Evaluation planning

**MUO-E3036 Interaction Design (IxD)**

**31 January 2022**

**Antti Salovaara**

MyCourses > Interaction design > Split S > Lecture slides >  
Week4-Day1-Evaluation-planning.pdf

# Contents of today's teaching

Choosing what features to compare in evaluation

Planning the evaluation

Lots of group work during the day:

- Discussion on which designs you want to compare

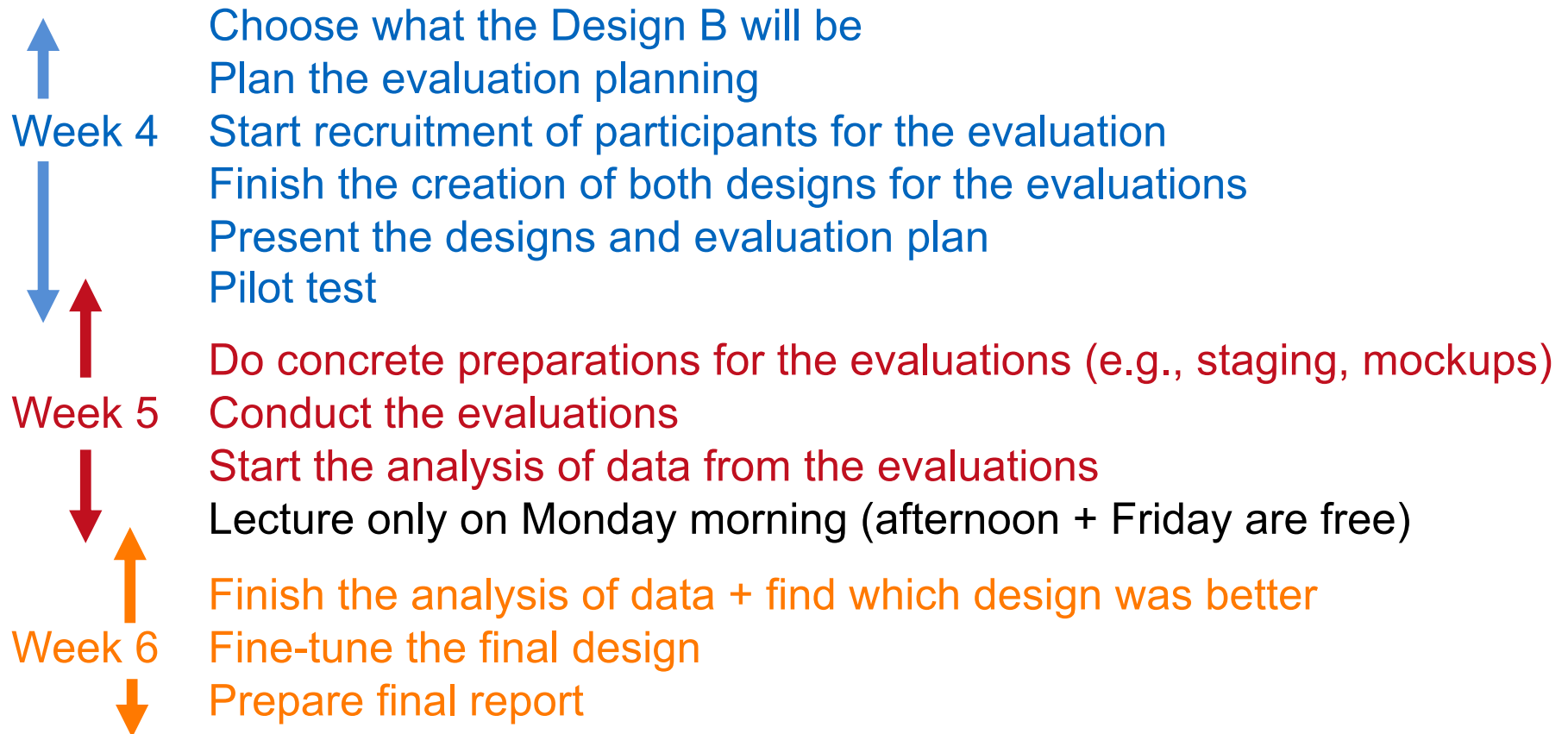
- Discussions on useful evaluation arrangements

Friday's presentation instructions

Reading materials for Friday

Tutor meetings

# Contents for all the remaining weeks



↑  
Fuzzy boundaries

# Friday's presentation instructions

# Friday's presentation contents

10 minutes / group

## Part 1: A vs B presentation

Your UX goal(s)

Final design A vs Final design B

How these designs address your UX goal in different ways

## Part 2: Evaluation plan

How you will evaluate the designs with users

How you will measure the UX goal(s)

Also: Submission of the evaluation plan to MyCourses

<https://mycourses.aalto.fi/mod/assign/view.php?id=861451>

DL: Friday 13:00

# Evaluation plan template

## 1. Your UX goal(s): ½ page

Name each goal + tell using your own words what it means in the case of your app or service

## 2. Present your designs (A and B): 1–2 pages

Screenshots from your final designs + main interactions within and between the screens

Clearly indicate + explain what makes A vs B different

Explain the reasons for the two designs (e.g., how the differences related to your UX goals)

## 3. Evaluation plan: 2–4 pages

Details of your methods: interview questions, usability test script (and division of work), questionnaires

How these methods will answer which design (A or B) meets your UX goal better

# Finishing the Designs A and B

Choosing what to compare

Designing A and B so that this comparison is possible

# Ways by which designs may differ

Different interaction sequences

“First step A, then B” vs “First step B, then A”

Different IxD patterns

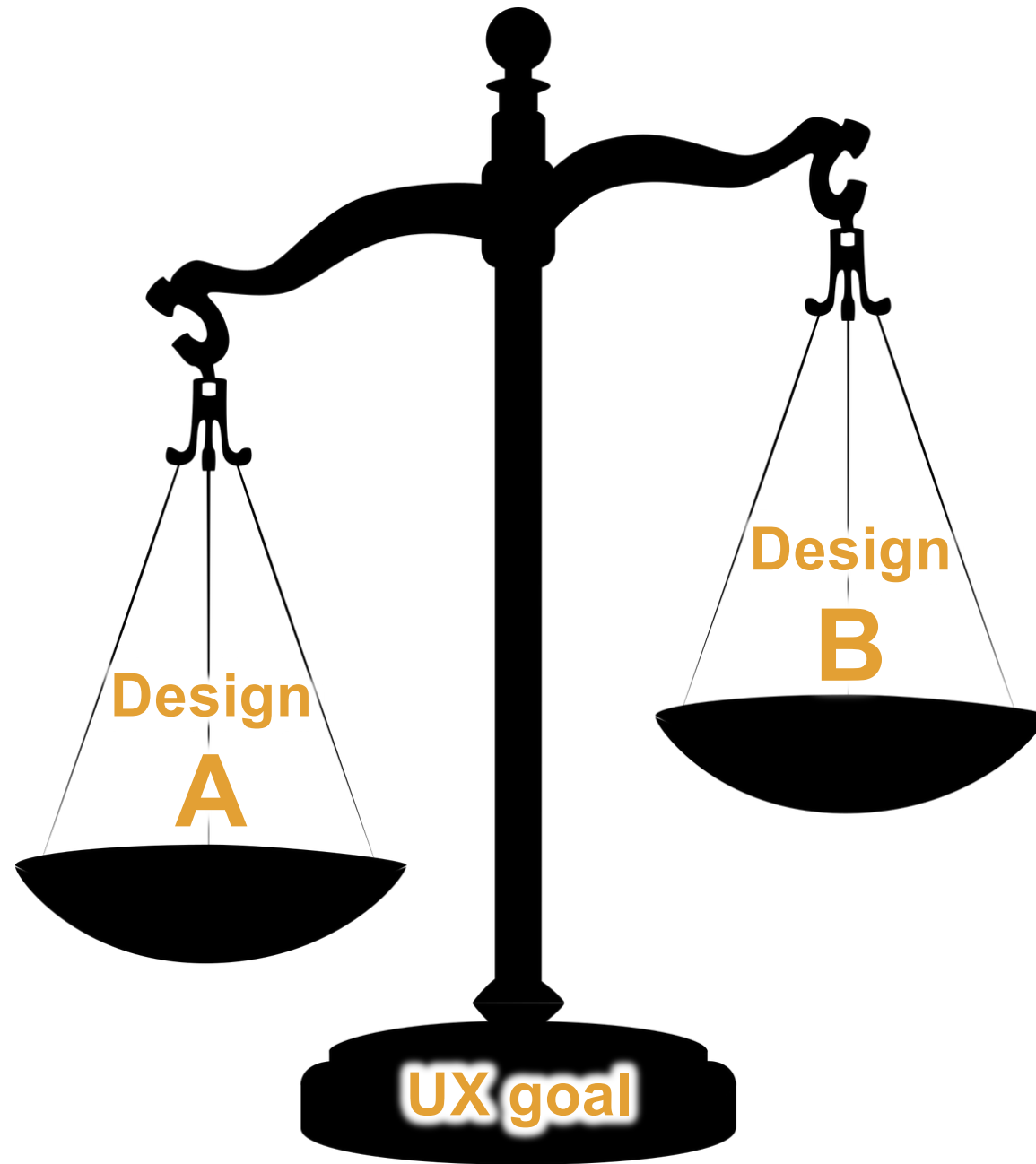
Wizard vs. accordion

Different solutions to the same problem

Different information visualizations

Etc.





# Discussion in groups (15 mins)

## Overall question:

“What kind of design B should we have?”

## Guiding questions:

What could be a better way to reach the UX goal than Design A?

What kind of difference would be worth of a user evaluation?

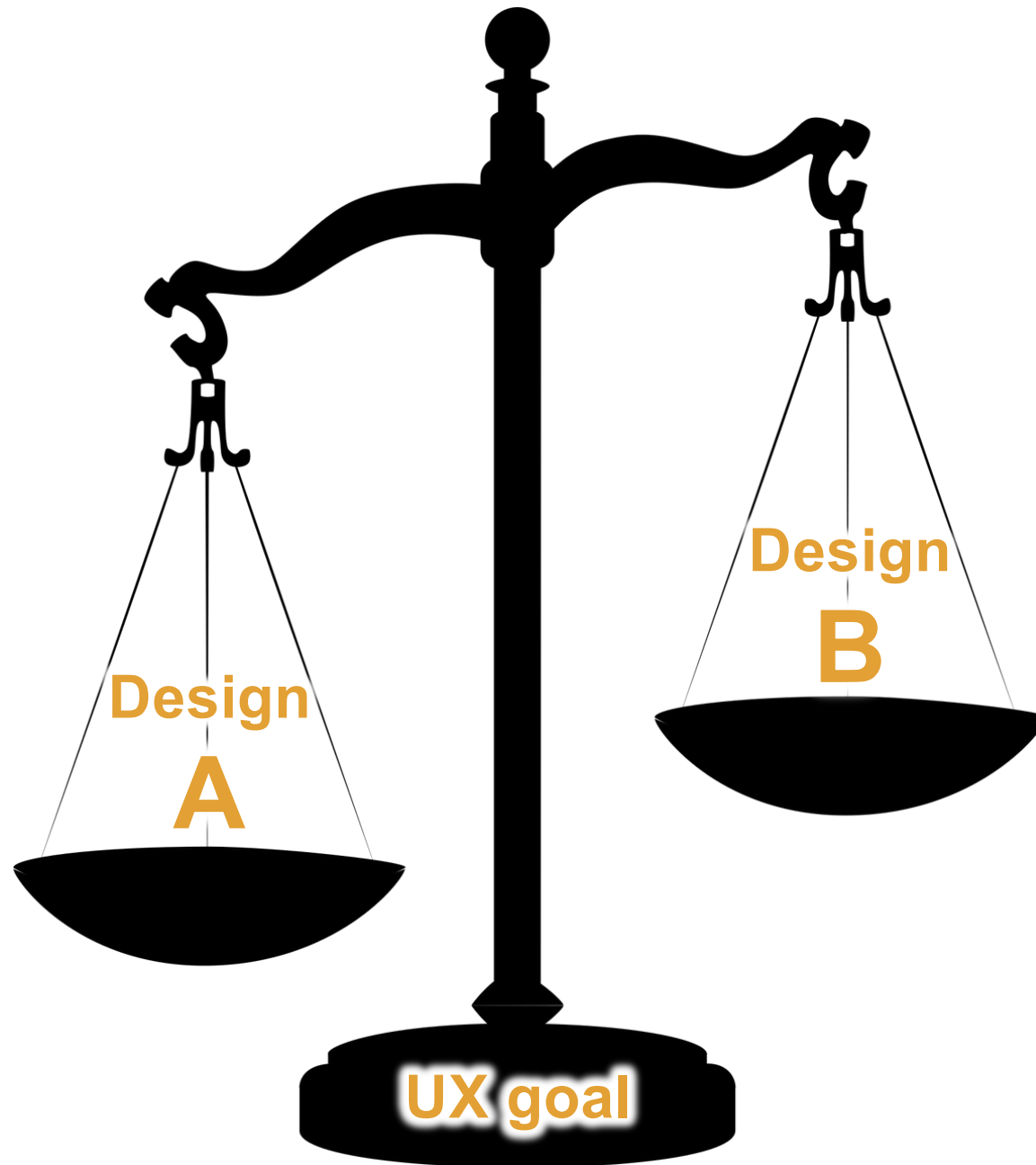
Join the main session if/when you wish to discuss with me

## After the discussion:

Groups present their ideas to each other (**11** ↔ 12, **13** ↔ 14, **15** ↔ 16)  
+ get feedback (15 mins)

Joint discussion about questions that emerged

# Planning the evaluation



What determines which design has more weight?

=> It is the data that you collect in your evaluation

# Quick group task

Consider your most important UX goal

E.g., ease of use

Discuss: What kinds of data do you need so that you can evaluate this UX goal?

E.g., user's stress level, number of errors at the first try on the task, ...

Steps:

1. Start by brainstorming individually (5 mins)
2. Then share ideas within your group (5 mins)

# Break

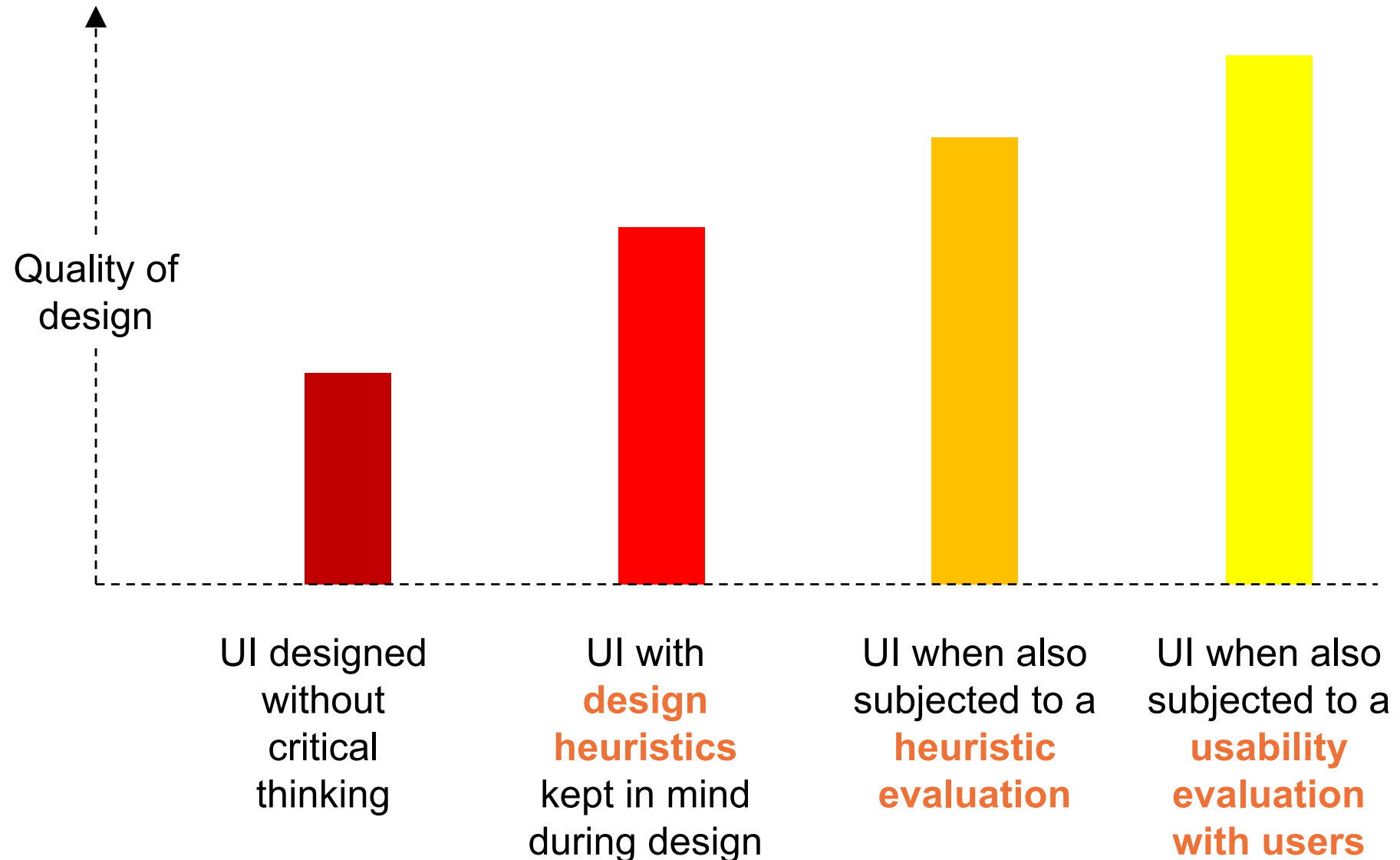
## Types of evaluations

Heuristic evaluation: evaluation without users

Traditional scenario-driven usability evaluation

In-the-wild evaluation

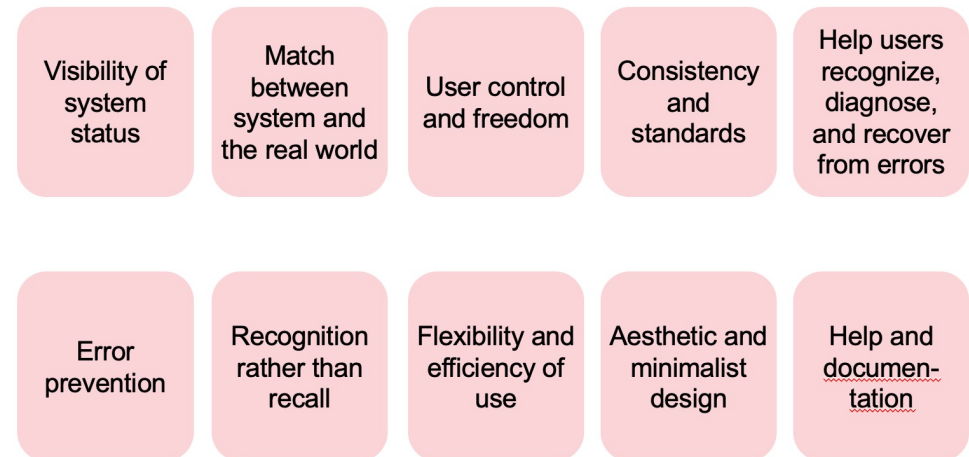
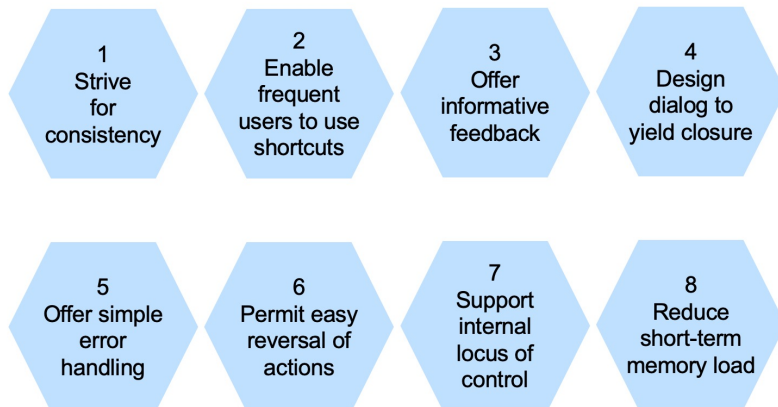
# How to reach a good usability and UX



# Heuristic evaluation

## UI's analysis using the design heuristics

- Use both knowledge in the world and in the head
- Simplify the structure of tasks
- Make things visible
- Get the mappings right
- Exploit the power of constraints
- Design for error
- When all else fails: Standardize!





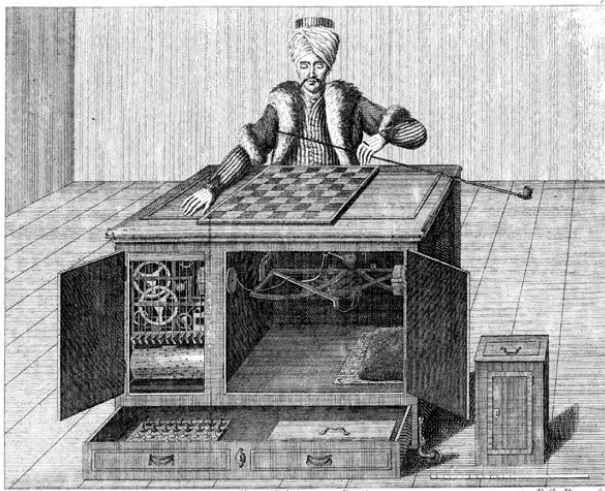
# Traditional usability evaluation



## Scenario-driven test:

1. Write realistic task scenarios for the features that need evaluation
2. Create mockup materials that make the unfinished system feel real
3. Present the scenario for the participant and ask him/her carry out the tasks.
4. Record with video
5. Repeat with more participants until findings “saturate”

# Wizard-of-Oz evaluations



Chess-playing automaton constructed by Wolfgang von Kempelen in 1770

Copper engraving in Karl Gottlieb von Windisch, Briefe über den Schachspieler des Hrn. von Kempelen, nebst drei Kupferstichen die diese berühmte Maschine vorstellen. 1783. Public domain: Copyright expired.

## Definition:

“a research experiment in which subjects interact with a computer system that subjects believe to be autonomous, but which is actually being operated or partially operated by an unseen human being.”

(Wikipedia)

## Use when:

you can't prototype a computer to perform interactions

## Ethics issue:

Setup is revealed after the study



<https://hcd498processlog.wordpress.com/2015/05/11/wizard-of-oz-a-pen-that-corrects-you-when-you-write-off-line/>

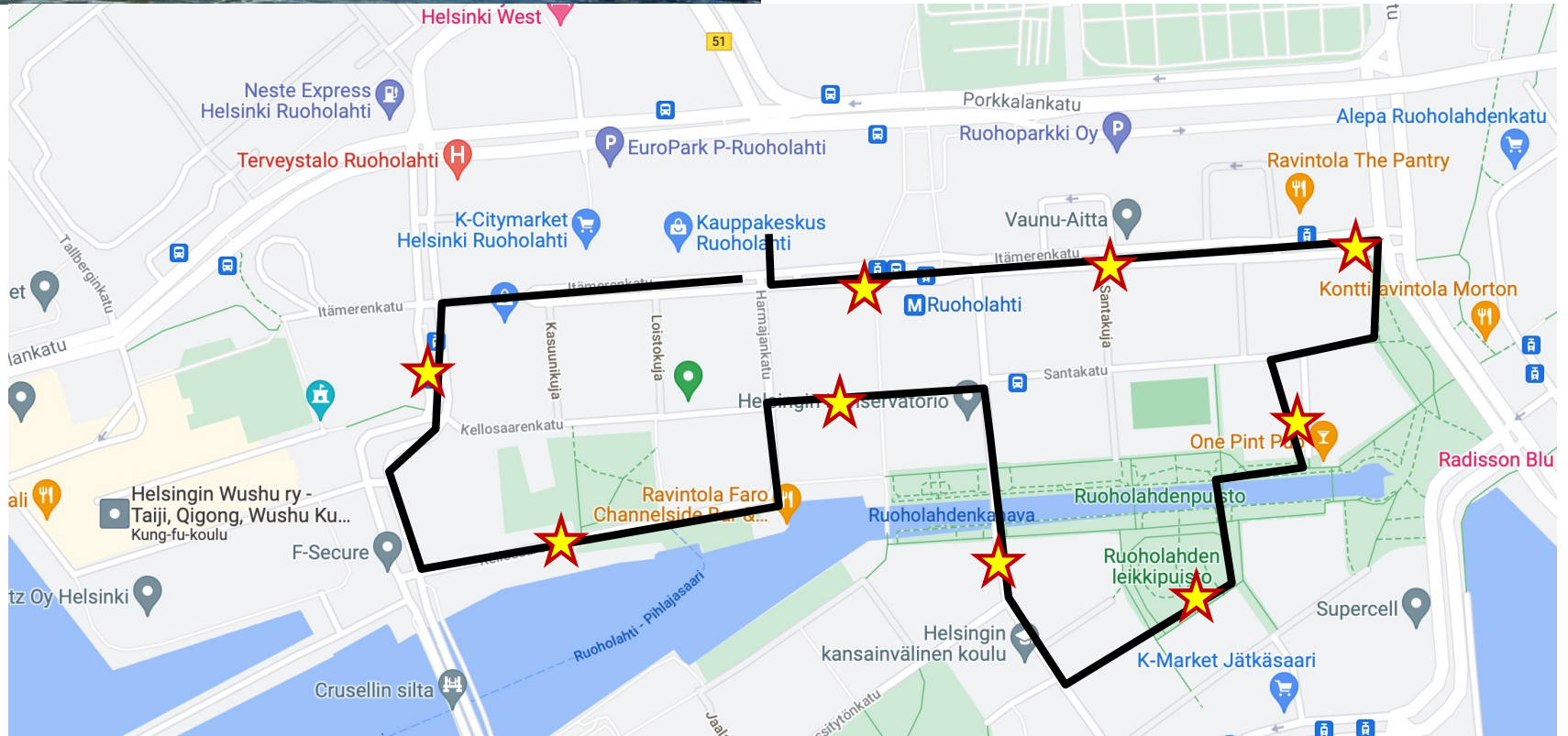
# Example: In-the-wild Wizard-of-Oz study

Would parents with babies be interested in location-based advertisements?

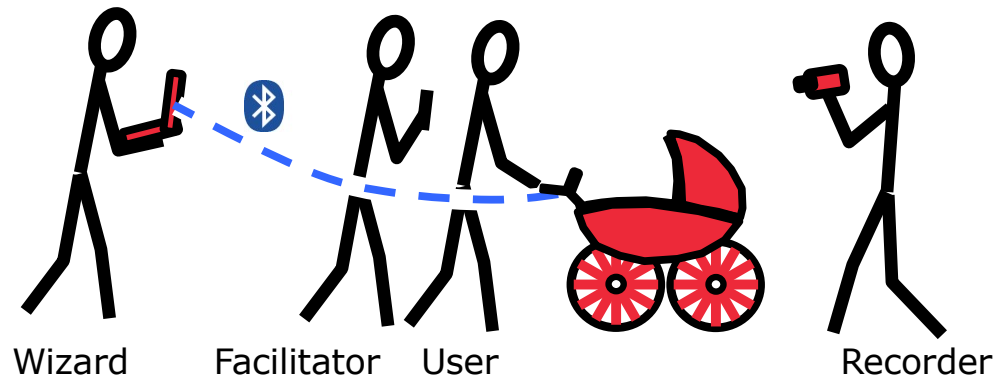




The Ruoholahti canal by Oghmoir.  
[http://commons.wikimedia.org/wiki/File:Ruoholahden\\_kanava.jpg](http://commons.wikimedia.org/wiki/File:Ruoholahden_kanava.jpg). Licensed under Creative Commons Attribution-Share Alike 3.0 Unported



# Wizard-of-Oz setup



## Wizard's controls



# How to control a prototype remotely

“In-the-wild” wizard-of-Oz studies on mobile phones:

It is a great benefit if you can make user’s screen contents change at desired times

Prepare a phone for the participant and allow its control from another phone

Investigate these options:

<https://joyofandroid.com/how-to-remotely-control-android-phone/>

<https://www.androidauthority.com/how-to-remote-control-android-device-41969/>

# Think aloud method

Origins in psychological research on problem-solving and creativity\*

Encourage the users to talk aloud:

- What they are trying to do

- What they are thinking

**!!!** Thinking aloud is not natural to many people

- A demonstration by the moderator and a practice task are needed to give the user an idea on what is expected

- Remember to remind the user politely (“Can you tell what you are now thinking?”)

\* E.g., Ericsson, K. A. (2006). Protocol analysis and expert thought: concurrent verbalizations of thinking during experts' performance on representative task. In K. A. Ericsson, N. Charness, P. J. Feltovich, & R. R. Hoffman (Eds.), *Cambridge Handbook of Expertise and Expert Performance*, ch. 13 (pp. 223--242). Cambridge University Press.

# Break?

## Arrangement of a user evaluation

Who should be recruited as users

How much data is needed



# What users should be recruited?

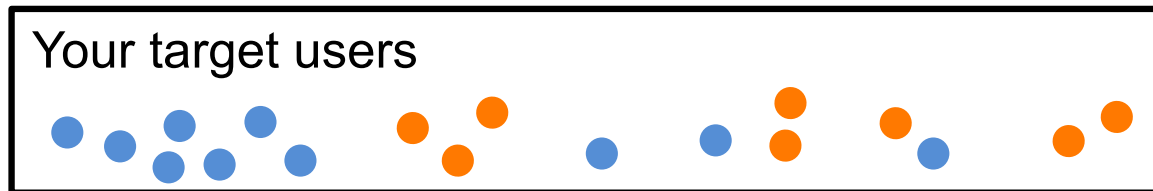
## Random sampling

Each participant that you recruit has a **known probability** of being chosen for the study

**Practically impossible** in studies on humans

## Convenience sampling

Studying people who you have a good access to (the typical method)



## Choosing between heterogeneous vs homogeneous samples

Homogeneous (users very similar): If you need “deep” findings

Heterogeneous (users differ a lot): Generalizable but shallower findings

# Choose between heterogeneous vs homogeneous samples

## Homogeneous sample:

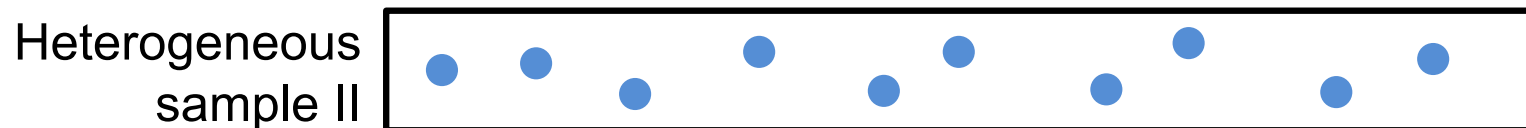
Users are very similar

Little noise in your data => You can get “deeper” findings

## Heterogeneous sample:

Users differ a lot (e.g., in terms of age, gender, expertise, life values)

A lot of noise and variability => Generalizable but shallower findings



# Using same participants again?

Pros and cons of using your Sprint's test participants again:

## Pros:

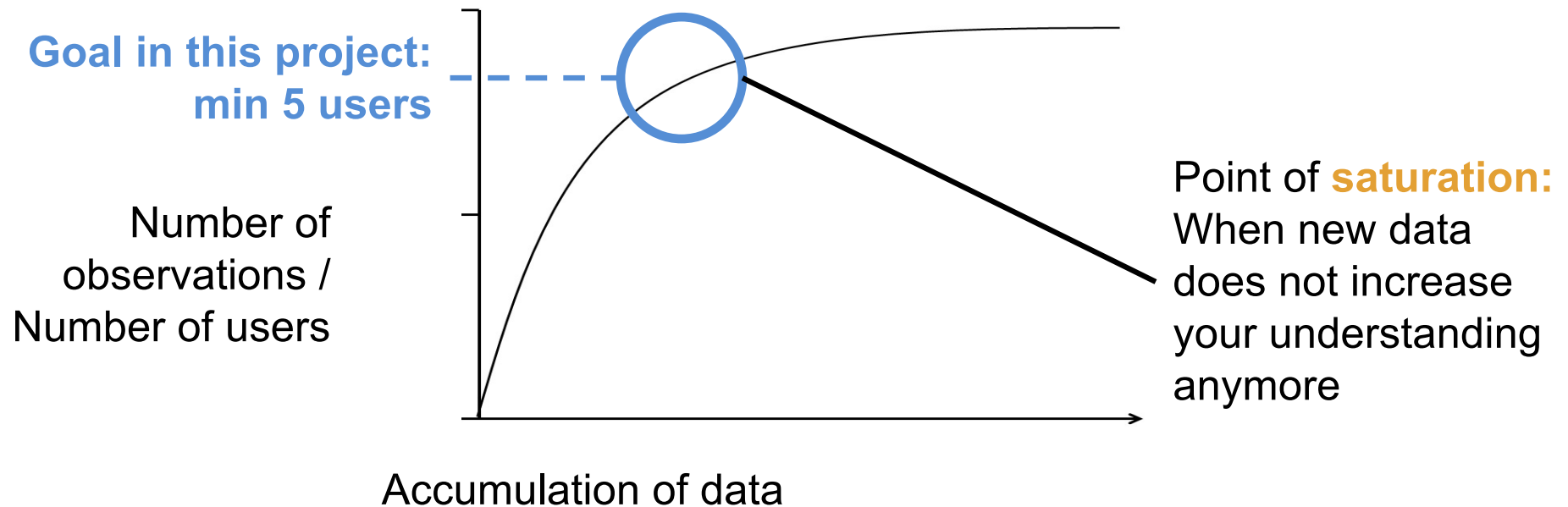
- More detailed feedback
- Easier recruitment
- No need to explain prototype in detail

## Cons:

- Overfitting your design to individual users' needs
- Learning effect
- May go against your UX goal evaluation (e.g., is it possible to evaluate ease of use with a user that already knows the product?)

Recruitment from this course: Same issues

# How much data is enough?



In usability evaluations, data can both quantitative and qualitative, but the analysis is almost always qualitative

# Making the most out of every participant

To gather more data, include repetition in the scenario

E.g., plan the first task to lead to suboptimal outcome, in order to make the user do something also another time

“Ok, now I have done almost what I wanted, but this is not perfect. I’ll try to find a better solution, just a minute...”

Gather data in many ways simultaneously:

Measure speed, errors etc.

Use think-aloud to also find out what the user thinks

Take video to observe behaviour and interactions

Use a questionnaire (SUS, AttrakDiff, your own questions...)

Interview about the experience after the evaluation

# Mockups and staging

Although evaluations are unnatural...

(since user are recruited to carry out artificially constructed tasks)

...they should feel natural and believable

(to help the participants engage in the tasks and behave naturally)

**Mockups: Preparation of authentic-feeling task materials**

=> To evaluate a CAD software, prepare an unfinished 3D design that the user can work on

**Staging: Making believable physical and social surroundings**

=> To evaluate a wayfinding app for busy shopping malls, you have to create a context of a busy shopping mall

# Comparative evaluations

A/B tests

Between-subjects vs within-subjects research designs

# A/B tests

## Gold standard for A/B tests:

A randomized statistical test between two systems that differ only by one factor

Example: in an online service, 10,000 visitors are directed to Design A, another set of 10,000 visitors are directed to Design B. The length of the visit is measured to find out which Design keeps them longer at the service.

Read more: Wikipedia: “A/B testing”, Google: “ab test ux design”

## In our course, the A/B test will be qualitative

What people say about Design A vs Design B

How they use the designs differently

Etc.

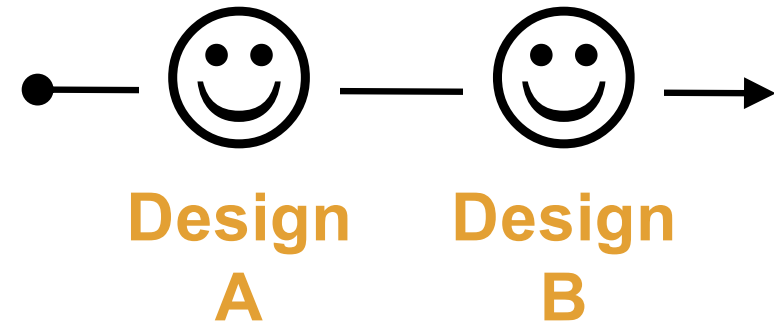


# What designs will each user interact with?

Participant uses both Design A and Design B



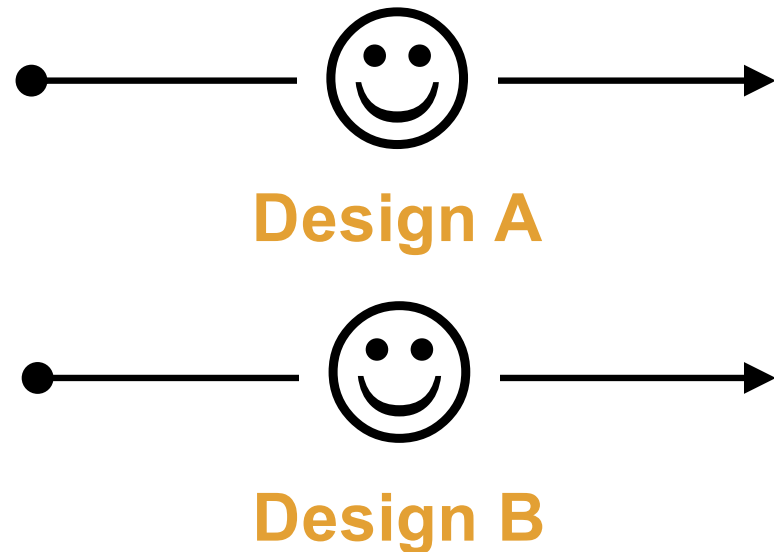
**Within-subjects**  
(aka repeated measures)



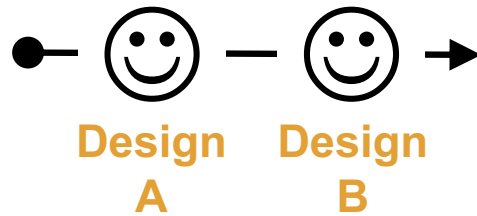
Participant user either Design A or Design B, but not both



**Between-subjects**

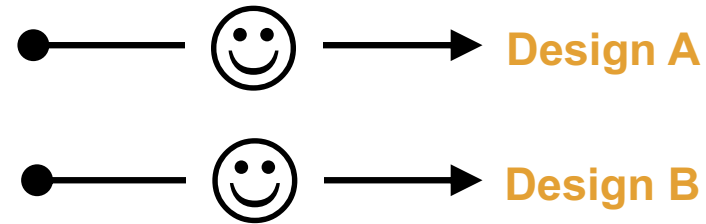


# Pros and cons of between- and within-subjects tests



## Within-subjects

- + A and B can be compared easily on user-by-user level
- + You get more data with a small number of people
- Learning effect: participants learn to carry out Task B by carrying out Task A



## Between-subjects

- + No learning effects
- Need more participants

← **Counter-balancing** helps:  
50% of users start with Design A,  
the other 50% with Design B

**Finally:**

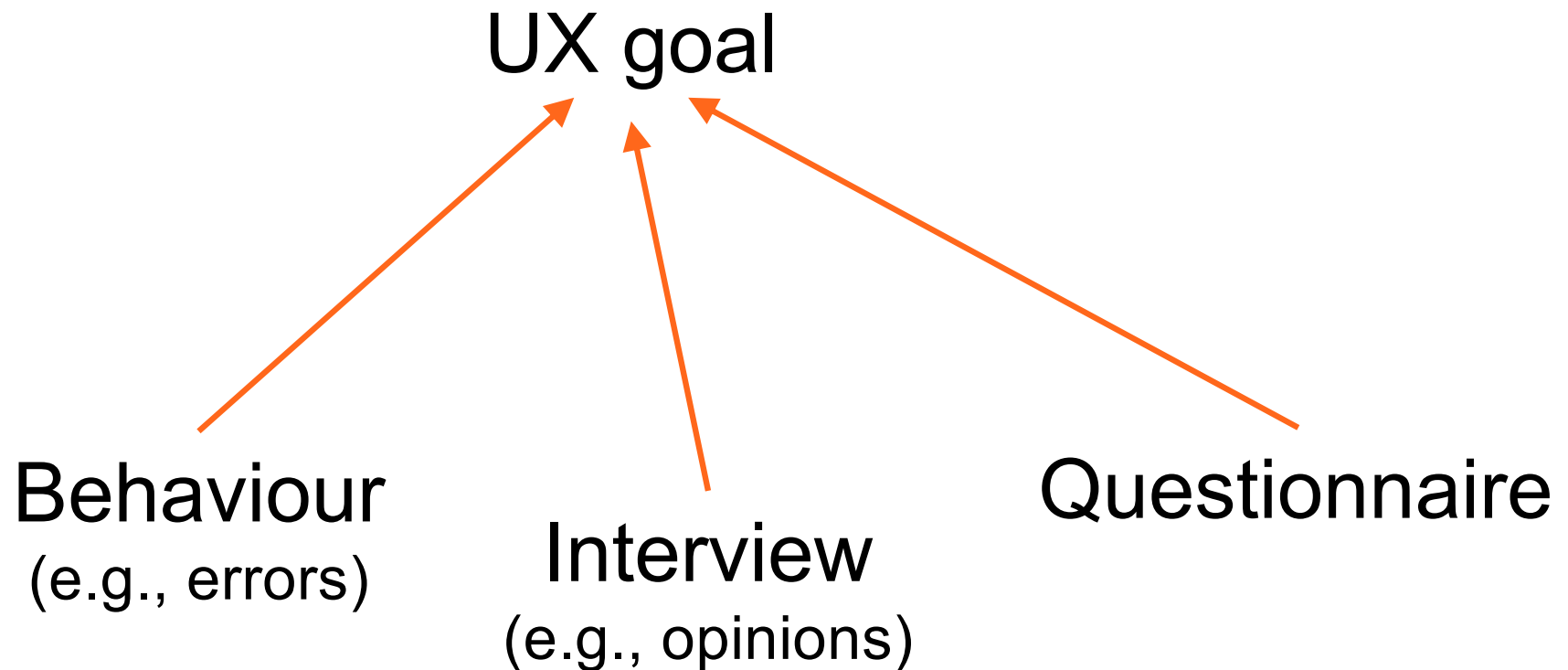
**What data should you gather?**

Triangulation

Examples of basic usability metrics

# Triangulation

Try to measure the same question in several, complementary ways



# Group your ideas from the quick group task

You did this task:

## Quick group task

Consider one of your UX goals

E.g., ease of use

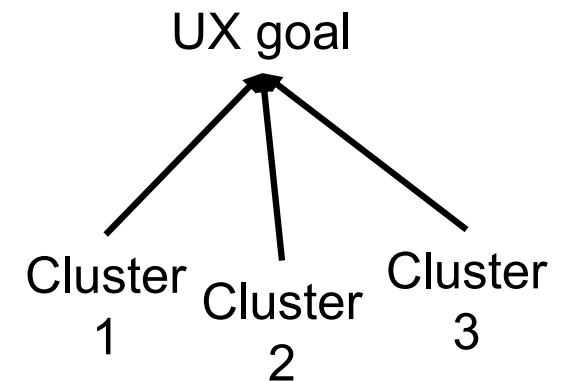
What kinds of data do you need so that you can evaluate this UX goal?

E.g., user's stress level, number of errors at the first try on the task, ...

Steps:

1. Start by brainstorming individually (5 mins)
2. Then share ideas within your group (5 mins)

Group your ideas into clusters that support each other



If you notice missed opportunities, add them to your lists

# Common usability metrics

**Table 10.3** Common usability metrics

Usability objective	Effectiveness measures	Efficiency measures	Satisfaction measures
<b>Overall usability</b>	<ul style="list-style-type: none"> <li>● Percentage of tasks successfully completed</li> <li>● Percentage of users successfully completing tasks</li> </ul>	<ul style="list-style-type: none"> <li>● Time to complete a task</li> <li>● Time spent on non-productive actions</li> </ul>	<ul style="list-style-type: none"> <li>● Rating scale for satisfaction</li> <li>● Frequency of use if this is voluntary (after system is implemented)</li> </ul>
<b>Meets needs of trained or experienced users</b>	<ul style="list-style-type: none"> <li>● Percentage of advanced tasks completed</li> <li>● Percentage of relevant functions used</li> </ul>	<ul style="list-style-type: none"> <li>● Time taken to complete tasks relative to minimum realistic time</li> </ul>	<ul style="list-style-type: none"> <li>● Rating scale for satisfaction with advanced features</li> </ul>
<b>Meets needs for walk up and use</b>	<ul style="list-style-type: none"> <li>● Percentage of tasks completed successfully at first attempt</li> </ul>	<ul style="list-style-type: none"> <li>● Time taken on first attempt to complete task</li> <li>● Time spent on help functions</li> </ul>	<ul style="list-style-type: none"> <li>● Rate of voluntary use (after system is implemented)</li> </ul>
<b>Meets needs for infrequent or intermittent use</b>	<ul style="list-style-type: none"> <li>● Percentage of tasks completed successfully after a specified period of non-use</li> </ul>	<ul style="list-style-type: none"> <li>● Time spent re-learning functions</li> <li>● Number of persistent errors</li> </ul>	<ul style="list-style-type: none"> <li>● Frequency of reuse (after system is implemented)</li> </ul>
<b>Learnability</b>	<ul style="list-style-type: none"> <li>● Number of functions learned</li> <li>● Percentage of users who manage to learn to a pre-specified criterion</li> </ul>	<ul style="list-style-type: none"> <li>● Time spent on help functions</li> <li>● Time to learn to criterion</li> </ul>	<ul style="list-style-type: none"> <li>● Rating scale for ease of learning</li> </ul>

Benyon  
p. 226

# Ask comparative questions

Use the comparative setup to gather deep answers:

“If you would need to analyse **how easy** these Designs were to use, how would you describe them?”

“Can you tell 2 good aspects from both Designs? How about 2 negative aspects?”

“How many stars, from 1 to 5, would you give to these Designs along the following dimensions: ease of use, efficiency, simplicity, beauty. Explain why.”

# Questionnaire-based measures

SUS

AttrakDiff

NASA TLX



# System Usability Scale (SUS)

## Usability.gov's description:

“Quick and dirty”, reliable tool for measuring the usability. It consists of a 10 item questionnaire with five response options for respondents; from Strongly agree to Strongly disagree.

<https://www.usability.gov/how-to-and-tools/methods/system-usability-scale.html>

## Example statements:

1. I think that I would like to use this system frequently.
3. I thought the system was easy to use.
6. I thought there was too much inconsistency in this system.

These can be great discussion topics after the user has given their responses

# AttrakDiff

<http://attrakdiff.de/index-en.html>

Contains a web tool to carry out all the analyses

AttrakDiff measures users' perceptions with 28 “semantic differentials”:

Ugly — Beautiful

Confusing — Clear

...

**Result: three measures:**

Pragmatic (utilitarian) quality

Hedonic (enjoyment-oriented) quality

Attractiveness

**Check out the use for A/B tests:**

<http://attrakdiff.de/index-en.html#tab-vergleich-ab>

# NASA TLX (task load index)

Measures subjective perception of task load

Traditional version:

6 statements

Ranking of the statements task

Score calculation

“Raw NASA”:

Plain average of the 6 statements

More info + where to get it:


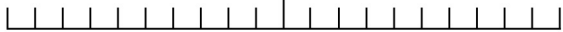
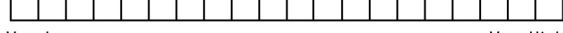
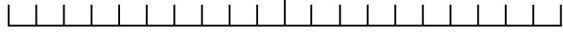
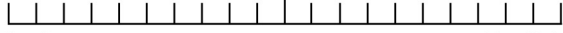
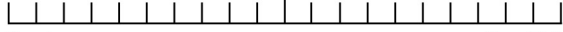
<https://humansystems.arc.nasa.gov/groups/TLX/>

<https://en.wikipedia.org/wiki/NASA-TLX>

Figure 8.6

## NASA Task Load Index

Hart and Staveland's NASA Task Load Index (TLX) method assesses work load on five 7-point scales. Increments of high, medium and low estimates for each point result in 21 gradations on the scales.

Name	Task	Date
Mental Demand	How mentally demanding was the task?	
Very Low		Very High
Physical Demand	How physically demanding was the task?	
Very Low		Very High
Temporal Demand	How hurried or rushed was the pace of the task?	
Very Low		Very High
Performance	How successful were you in accomplishing what you were asked to do?	
Perfect		Failure
Effort	How hard did you have to work to accomplish your level of performance?	
Very Low		Very High
Frustration	How insecure, discouraged, irritated, stressed, and annoyed were you?	
Very Low		Very High

# Your own questionnaire

If you want, you can make your own questionnaire

Tip:

- Use Likert statements (“Totally disagree – Totally agree”)

- Ask about the same topic using multiple prompts

- Do a pilot study

# How to visualize questionnaire data

Use bar charts to visualize answers

1. Calculate user-level averages:

Example: Average of all NASA-TLX answers from a user to Design A, and another average from answers to Design B

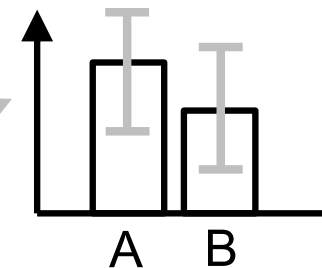
2. Calculate averages across all the users

Example: Average of all the user-level averages from step 1 for Design A, and the similar average for Design B

3. Present the two bars in a diagram

4. Draw confidence intervals

If you have time – See YouTube tutorials



# Group task: “Methods shopping”

Review the slides in the group

Pick measures that you think you would like to have

Discuss if this is possible

Discuss why the measure would be useful

Does it relate to your UX goal(s)?

Is it important because of another reason?

# How to “survive” this week

# Divide yourselves into sub-teams

Monday

Tuesday

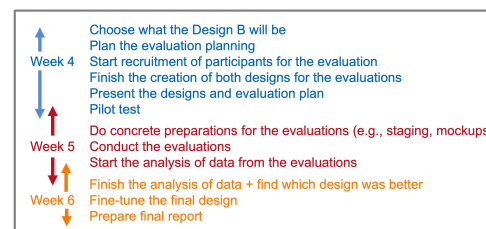
Wednesday

Thursday

Friday



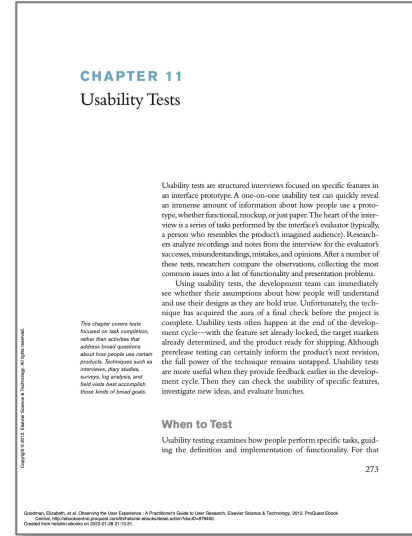
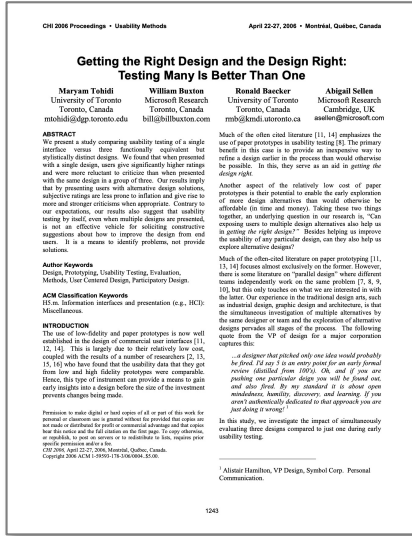
Remember that weeks 4, 5 and 6 also have fuzzy boundaries





# Reading materials

# Reading materials for week 4



Tohidi et al (CHI2006):

Getting the right design and the design right: Testing many is better than one  
<https://dl.acm.org.libproxy.aalto.fi/doi/10.1145/1124772.1124960>

Goodman & Kuniavsky (2012):

Chapter 11: Usability tests  
<https://pdfroom.com/books/observing-the-user-experience-second-edition-a-practitioners-guide-to-user-research/wW5mwke4gYo>  
or  
[https://primo.aalto.fi/permalink/358AALTO\\_INST/ha1cg5/alma998568944406526](https://primo.aalto.fi/permalink/358AALTO_INST/ha1cg5/alma998568944406526)

# Tutor meetings

[https://doodle.com/poll/xvr4fgimhs6w8m9f?utm\\_source=poll&utm\\_medium=link](https://doodle.com/poll/xvr4fgimhs6w8m9f?utm_source=poll&utm_medium=link)