



# EMPIRICAL USER STUDIES

---

*Evaluation of the Anima tool*  
*Rui Couto • José C. Campos*

# OUTLINE

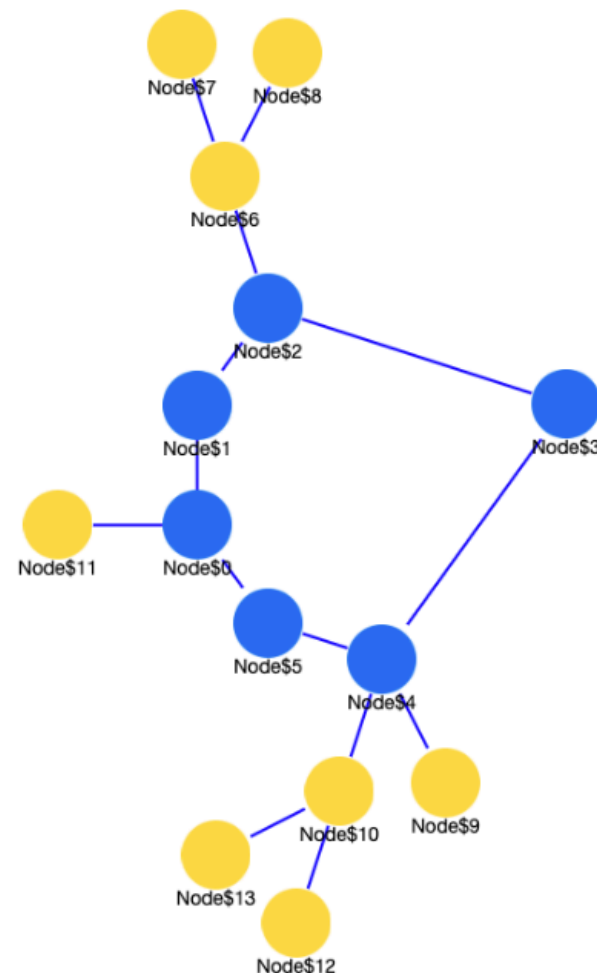
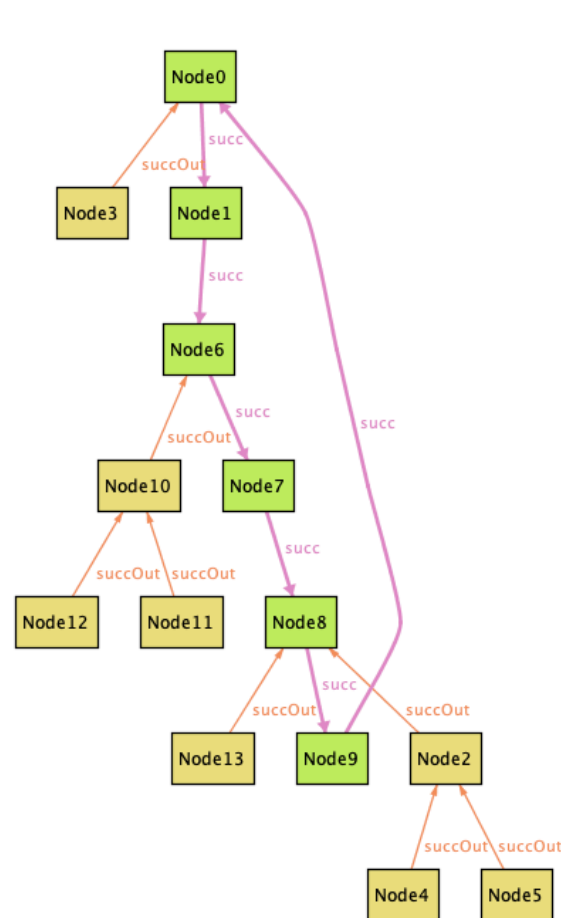
---

- The problem
- Objective specification
- Study design
- Elaboration
- Evaluation
- Publication
- Case study: OutSystems Learnability model

# THE PROBLEM

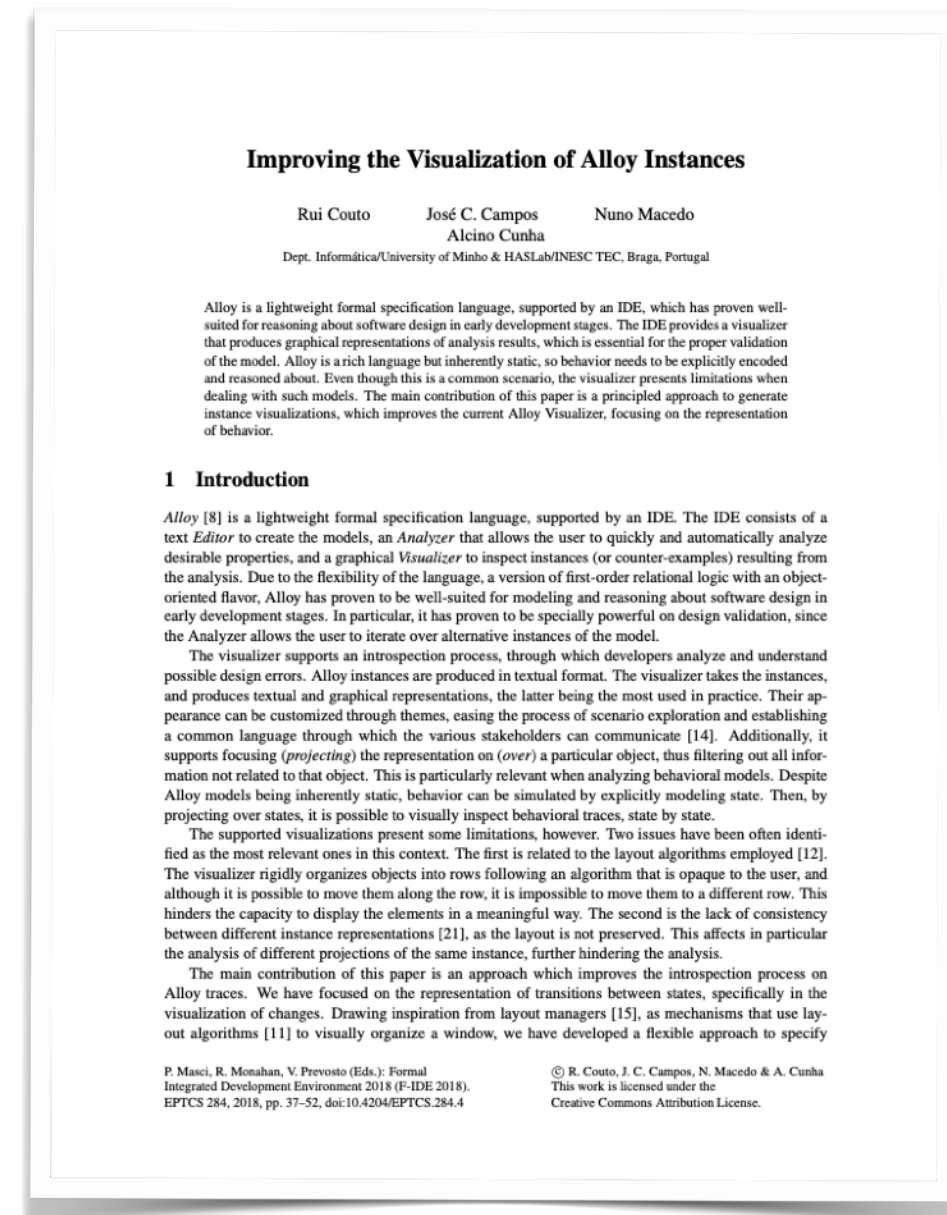
.....

- The motivation behind a user study is an existing problem in a given software product
- In our case, we have a concrete problem:
  - The Alloy visualiser is hard to use, not appealing, and we believe it leads to errors.
- Our solution was to develop Anima, a new visualiser which improves the Alloy visualiser, with features to tackle the identified issues



# THE PROBLEM

- Since Anima was designed to solve the Alloy visualiser issues, we believe to have a better solution<sup>1</sup>.
- However...
  - We don't know for sure if such is true, or if it is only our perception - we have developed the tool
  - A small group of users, which know the problems in beforehand isn't a significative group
  - Solving the identified issues is not a guarantee to have a better solution: introduction of new problems, solving false problems, etc.

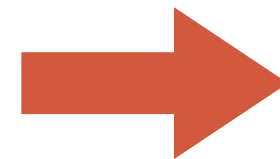
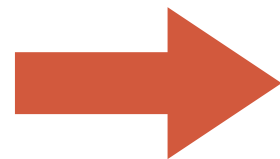


<sup>1</sup> Couto, Rui, et al. "Improving the Visualization of Alloy Instances." *arXiv preprint arXiv:1811.10817* (2018).

# THE PROBLEM

---

- The solution to address this problem is to:
  - Perform a user study with independent participants (reducing bias)
  - Measure different dimensions such as performance, time, and effort in both visualisers (gather objective data)
  - Compare the results (achieve concrete results)





# OBJECTIVES


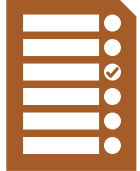


---

- The user study is shaped according to the objectives to achieve
- The objectives define the dimensions to be measured
  - E.g. How to evaluate a compiler? And a visualisation tool?
  - Are we measuring single performance, or comparing tools?
- In this case, we are evaluating if Anima was able to improve the analysis process, when compare with Alloy.

# OBJECTIVES

---

- We define that, in order for a visualiser to be better than the other, the following aspects should be considered:

-  **Time** - required to perform tasks
-  **Number of errors** - resulting from interpretation mistakes
-  **Effort** - required to perform the tasks
-  **Overall perception** - that the users had for both tools

- Having shaped the objectives, we were able to design the study accordingly.

# OBJECTIVES

---



## Time

- If a user takes less time to perform a task, then the tool is more effective
- More time spent in a task, might result in fatigue, and increase the number of errors



## Number of errors

- These tools have the objective of showing possible errors in systems
- Performing an error means the possibility for an error in the final solution.
- Having fewer errors is one of the most relevant objectives



# OBJECTIVES

---



## Effort

- A higher effort is associated with a higher probability for interpretation errors
- A lower effort in the interpretation means the user can focus in the interpretation of the problem itself



## Overall perception

- While not an easily quantifiable measure, if the users have a bad perception of a tool, they will be less prone to use it
- A higher perception might result in a higher adoption ratio

# STUDY DESIGN

---

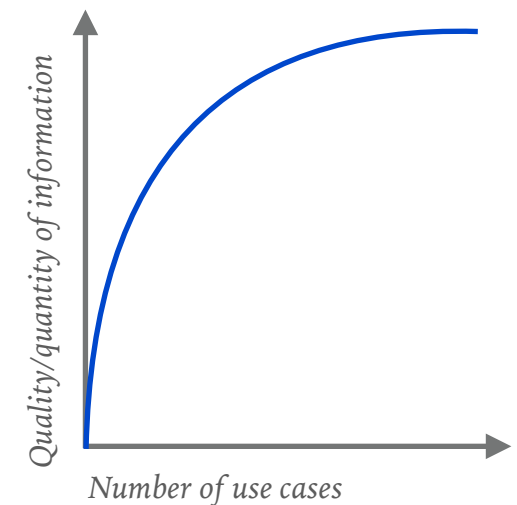
- Designing the study consists in:
  1. Defining the case study to be performed
  2. Defining how to measure the defined dimensions
  3. Predict possible problems
  4. Performing pilot studies



# STUDY DESIGN – 1 DEFINING THE CASE STUDY TO BE PERFORMED

.....

- The case study should be relevant, i.e. illustrate real usage scenarios
- If we design a use case based in the weaknesses/strength of our tool, we will be hindering/promote results of the other tool
- Higher number of studies provides more uniform results
  - Time to perform a study is limited, after  $\sim 1.5$ h users will start to show signs of fatigue, and results will not be useful
- Lower number of results provides less objective results
  - Results from two use cases can be drastically different



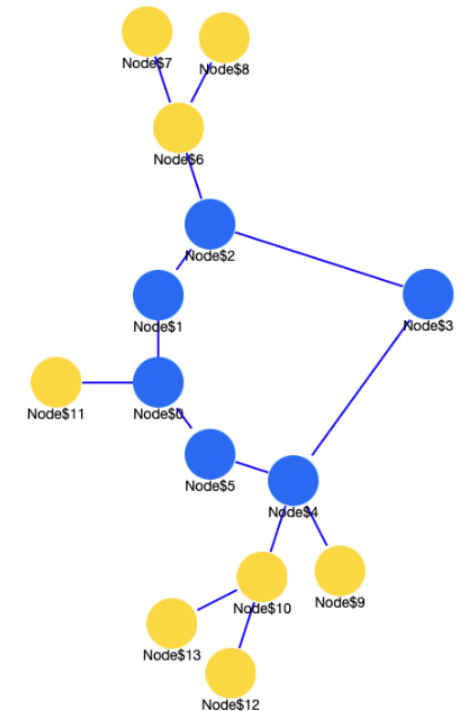
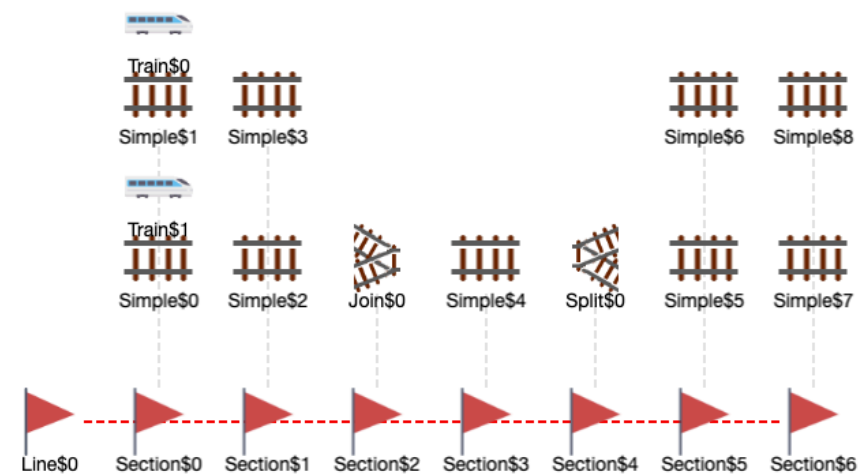
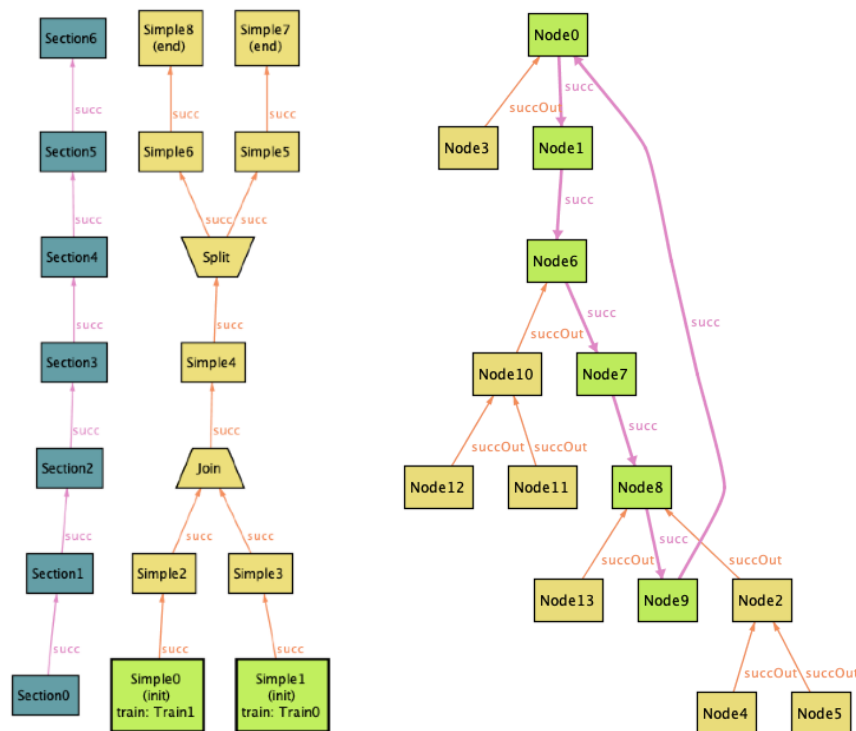
# STUDY DESIGN – 1 DEFINING THE CASE STUDY TO BE PERFORMED

---

- We have used two use cases, different in nature and context
  - A higher number of use cases would be too costly
  - Fewer user cases would provide biased results towards the use case
- We required, at least two scenarios, in order to test the customisation capabilities of Anima.

# STUDY DESIGN – 1 DEFINING THE CASE STUDY TO BE PERFORMED

- We have used two existing scenarios:
  - ERTMS/ETCS level 3<sup>1</sup> - a model of a train line, previously published
  - Chord<sup>2</sup> - A well known model in the Alloy community

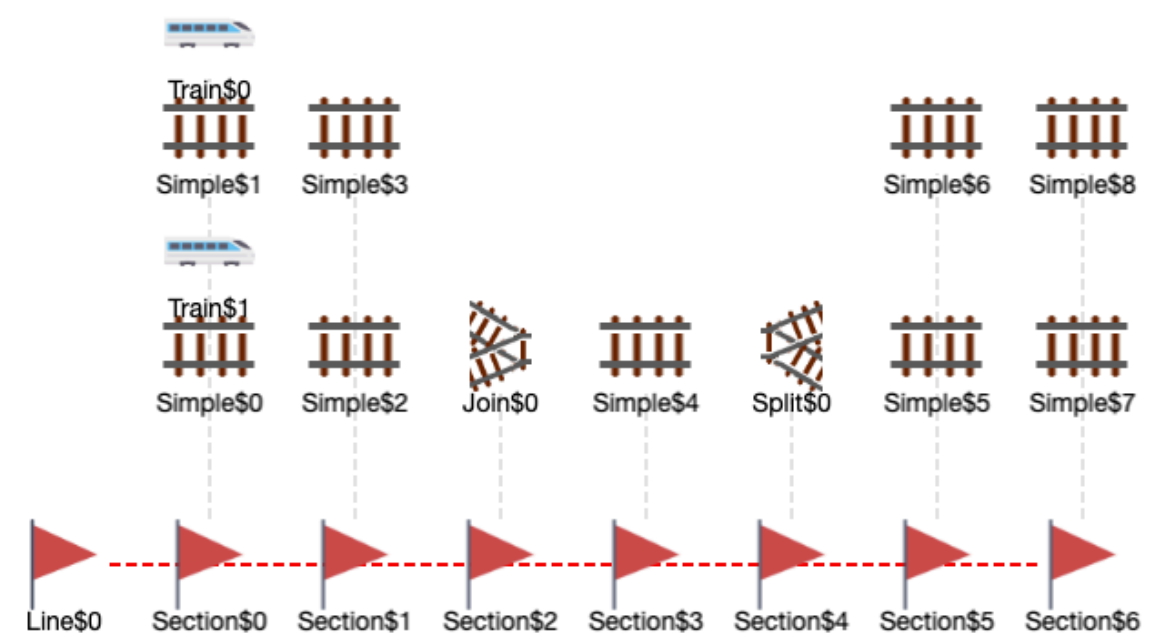
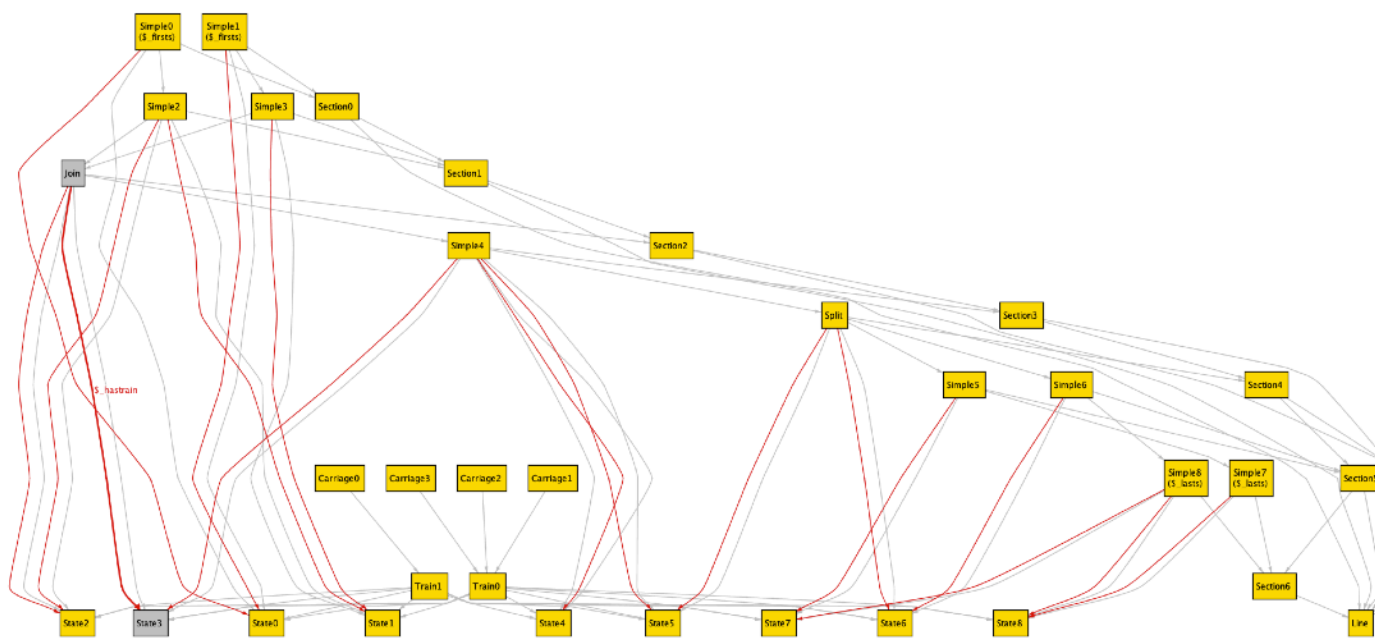


1. A. Cunha and N. Macedo. Validating the hybrid ERTMS/ETCS level 3 concept with electrum. In M. Butler, A. Raschke, T. S. Hoang, and K. Reichl, editors, *Abstract State Machines, Alloy, B, TLA, VDM, and Z*, pages 307–321, Cham, 2018. Springer International Publishing.

2. y, "Using lightweight modeling to understand Chord" (Pamela Zave; ACM SIGCOMM Computer Communication Review, 42(2):50-57, April 2012)

# STUDY DESIGN – 1 DEFINING THE CASE STUDY TO BE PERFORMED

- Fairness should be considered as well.
- Both Anima and Alloy tools allow for customisation
  - We should not use a non-themed version of the Alloy representation, vs a themed version of Anima
  - We have asked an Alloy expert to create the former one



1. A. Cunha and N. Macedo. Validating the hybrid ERTMS/ETCS level 3 concept with electrum. In M. Butler, A. Raschke, T. S. Hoang, and K. Reichl, editors, *Abstract State Machines, Alloy, B, TLA, VDM, and Z*, pages 307–321, Cham, 2018. Springer International Publishing.

2. y, "Using lightweight modeling to understand Chord" (Pamela Zave; ACM SIGCOMM Computer Communication Review, 42(2):50-57, April 2012)



# STUDY DESIGN – 2 DEFINING HOW TO MEASURE THE DEFINED DIMENSIONS

.....

- In order to produce relevant information, different outputs should be gathered.
- Some are objective:
  - Time is “easily” measured, as it can be objectively observed
  - Registering the number of errors is also an objective process
- We have measured time and errors as directly observable outputs
- Other are harder to measure:
  - Is the user understanding the problem?
  - Is the user under stress?
- We have observed the users to take further conclusions.

## STUDY DESIGN – 2 DEFINING HOW TO MEASURE THE DEFINED DIMENSIONS

.....

- Different techniques should be used.
- Direct observation is a non-intrusive, objective data gathering technique, e.g. time
  - For non objective data, interpretation is required, e.g. is the user understanding the problem?
- Interaction with the user helps to understand hard to observe aspects
  - Is the user really understanding the problem? Is the user struggling? Did the user make a mistake due to an interpretation error?
  - Risk of influencing the user, and affecting the final results.

## STUDY DESIGN – 2 DEFINING HOW TO MEASURE THE DEFINED DIMENSIONS

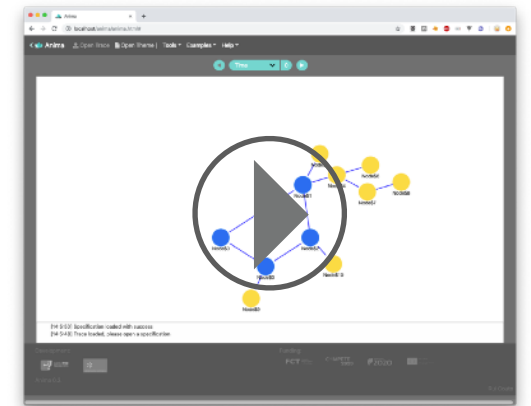
.....

- Standard questionnaires are well accepted
- They have been widely used, so a large amount of comparable information is available
- We decided to use two different questionnaires:
  - NASA TLX - measures the task load, i.e. how hard were the tasks to perform.
  - User Experience Questionnaire (UEQ) - measures the quality of the representations.
- These standard questionnaires allow us to disseminate our results in a way that other researchers understand

# STUDY DESIGN – 2 DEFINING HOW TO MEASURE THE DEFINED DIMENSIONS

.....

- Screen, audio and video are good sources for later analysis
  - They allow a deeper analysis, record statements e.g. “I like this feature”, or “This is hard to perform”
  - Can be viewed as many times as needed
- However...
  - They cause stress in the participants. Knowing that their interactions are being record, might affect their behaviour
  - Certain aspects (e.g. frustration) are easier to capture face to face
  - There are privacy concerns to be taken in consideration - some users don't like recordings, specially video; there are legal aspects to be taken in consideration



## STUDY DESIGN – 2 DEFINING HOW TO MEASURE THE DEFINED DIMENSIONS

---

- We tried to balance feedback record
  - We didn't record video, as we think it is the most intrusive
  - We record audio, as the participants did not seem affected by it
  - We recorded the screen, as it is the least intrusive recording media.
- Through direct observation and data logging we record remaining outputs.

# STUDY DESIGN – 3 PREDICT POSSIBLE PROBLEMS

---

- How to gather participants?
  - Consider the study context - we asked Interactive Systems students
- Is the study adequate for the target audience?
  - Anima is used for analysis in formal methods
- Is the time adequate?
  - If the study took too long, the participants are less willing to participate
- What if the users fail to perform the tasks?
  - We decided to guide participants in case of difficulties, taking care to avoid interfering with the results



# STUDY DESIGN – 3 PREDICT POSSIBLE PROBLEMS

---

- What are the learning effects?
  - We have created several combinations for the participants to start with different scenarios/tools
- What are the **threats to validity**?
  - Is the user focused? Is the user cooperating? (*Too much?!)*  
Some are out of our control
- How privacy problems can affect the results?
  - We clearly explained the context of the study, while respecting the users' privacy

- There is a set of documents required to perform the study.
- We must record consents, loggings, and results
- We should also provide visual aides to guide the participants
- In this study, we have prepared the following material:
  - For the observer:
    - A script to guide the study
    - A logger sheet - to record time and other outputs
  - For each user:
    - A privacy consent - to be signed
    - A form - for the users to provide answers
    - The UEQ

## Process

1. Read the context
2. Provide one of the following
  - a. Train\_1 in Anima; Train\_2 in Alloy + Chord\_1 in Anima; Chord\_2 in Alloy
  - b. Train\_1 in Anima; Train\_2 in Alloy + Chord\_2 in Anima; Chord\_1 in Alloy
  - c. Tia
  - d. Tia
3. Provide id
4. Provide UE
5. Provide TL

## Disclaimer

During this study it  
be published. The  
results will be app

## Script

I would like to star  
where you will ente

In the context of test  
analysis of errors  
the traces represent  
each trace.

You will be asked  
want it is worth i  
tools. (i.e., if you e  
how the tool could

Prior to the test, it  
For each scenario  
visualizer to analy  
during the study a

At the end and in  
Once again, we a  
performance.

Is there any conce

**Start by showing**

### Config A - Participant

Date: \_\_\_\_\_

#### Train\_1 - Anima

	start	end
P1	start	end
P2	start	end
P3	start	end
P4	start	end

#### Train\_2 - Alloy

	start	end
P1	start	end
P2	start	end
P3	start	end
P4	start	end

#### Chord\_1 - Anima

	start	end
P1	start	end
P2	start	end
P3	start	end

#### Chord\_2 - Alloy

	start	end
P1	start	end
P2	start	end
P3	start	end

[illegible]

# STUDY DESIGN – 4 PERFORMING PILOT STUDIES

---

- Performing a study right away is dangerous:
  - What if it fails to fit in the predicted time?
  - What if there are errors in the scripts?
  - What if we can clearly see that our tool does not improve over Alloy?



# STUDY DESIGN – 4 PERFORMING PILOT STUDIES

---



- Some studies are not repeatable - If we fail to perform them the first time, they cannot be repeated in the same context:
  - Users will be less willing to repeat the same study.
  - Users will learn about the tools, and know what to expect the next time.
  - A large amount of time is required to perform the studies!
- Pilot studies are a good approach to mitigate these problems
  - Measure the real expected time to perform the tasks
  - Identify possible errors and unexpected outcomes
  - Gather feedback (can even be considered in the final data)

# ELABORATION

---

- Performing the studies is a laborious task, which starts with the preparation
- All material should be printed and organised beforehand - all the time that can be saved in the study should be saved
- All the required material should be ready, namely the laptop, pen, mouse, software, room, etc. Even missing a mouse could seriously compromise a study, as the user might not be comfortable with the touchpad.
- A flexible schedule is required, as the users preferences should be the priority - users are more willing to participate.



# ELABORATION

.....

- The setup should be established prior to the test

*Laptop with Anima, Alloy and Active Presenter*

*Observer material*

*Tools for the user*



*User material*



# ELABORATION

---

- Three main steps occur in the elaboration:
  1. The context and details should be explained to the participant
  2. The participant performs the tasks, during which the observer will be taking notes
  3. The observer ends the study with a debriefing

# ANALYSIS

---

- After the first study, the results should be immediately analysed if possible
- The observer will remember easily the experiment
  - E.g. emotions, expressions
- Sometimes it is hard to take note of all feedback, but can be remembered afterwards
- If a mistake occurs, it is easier to fix (e.g. wrong annotation)

# ANALYSIS

- A detailed analysis process follows the study
- All the information should be represented in a proper way (mainly, spreadsheets)

	Version	Day	Hour	Age	Gender	Experience FM	Yrs Alloy	Most Attractive	Most Relevant	Which would use	Recommend
S01	A	T	12:00	21	M	0	0	Anima	Anima	Anima	yes
S02	B	T	13:30	21	F	0	0	Anima	Anima	Anima	yes
S03	C	T	15:00	22	M	0	0	Anima	Anima	Anima	yes
S04	D	W	09:00	22	M	0	0	Anima	Anima	Anima	yes
S05	A	W	09:30	21	M	0	0	Anima	same	Anima	yes
S06	B	W	10:00	21	F	0	0	Anima	Anima	Anima	yes
S07	C	W	10:30	27	M	0	0	Anima	Anima	Anima	yes
S08	C	W	11:00	22	F	0	0	Anima	Anima	Anima	yes
S09	D	W	12:00	22	F	0	0	Anima	Anima	Anima	yes
S10	A	W	13:30	22	M	0	0	Anima	depends	depends	yes
S11	B	W	14:00	21	M	0	0	Anima	Anima	Anima	yes
S12	D	W		21	M	0	0	depends	Alloy	depends	yes
S13	A	W		21	F	0	0	Anima	Anima	Anima	yes
S14	B	W		22	F	0	0	depends	Alloy	Alloy	no
S15	C	W		21	M	0	0	Anima	Anima	Anima	yes
S16	D	W		21	M	0	0	Anima	Alloy	Anima	yes
S17	A	W		21	F	0	0	Anima	Alloy	Anima	yes
S18	B	T	11:00	33	M	0	0	Anima	Anima	Anima	yes
S19	C	W	14:00	21	M	0	0	Anima	Anima	Anima	yes

- A detailed analysis process follows the study
- All the information should be represented in a proper way (mainly, spreadsheets)

	Version	Day	Hour	Age	Gender	Experience FM	Yrs Alloy	Most Attractive	Most Relevant	Which would use Recommend									
S01		Anima				Alloy				Anima				Alloy					
S02		Task 01				Task 01				Task 02				Task 02					
S03		Total	Total	Total	Total		Total	Total	Total		Total	Total	Total		Total	Total	Total		
S04		P1	P2	P3	P4	SUM	P1	P2	P3	P4	SUM	P1	P2	P3	SUM	P1	P2	P3	SUM
S06	S01	22	16	7	10	55	47	19	22	21	109	49	10	8	67	139	16	10	165
S07	S02	30	28	34	86	178	54	29	53	72	208	25	11	2	38	196	28	26	250
S08	S03	17	39	28	77	161	101	75	95	184	455	72	39	9	120	206	30	3	239
S09	S04	38	17	73	59	187	103	52	149	105	409	29	15	8	52	310	41	11	362
S10	S05	65	136	82	33	316	41	65	39	108	253	72	21	8	101	208	20	18	246
S11	S06	52	32	50	38	172	74	100	64	29	267	37	3	2	42	170	35	54	259
S12	S07	103	10	37	99	249	80	30	128	151	389	126	38	17	181	304	27	14	345
S13	S08	48	17	17	8	90	63	93	126	31	313	60	16	2	78	63	12	15	90
S14	S09	14	28	22	32	96	34	39	65	63	201	40	10	4	54	48	53	13	114
S15	S10	53	75	63	62	253	21	24	8	27	80	78	10	7	95	20	17	20	57
S16	S11	32	21	42	93	188	17	44	35	33	129	65	24	2	91	118	56	8	182
S17	S12	14	32	37	80	163	31	33	38	66	168	58	9	8	75	147	56	7	210
S18	S13	34	10	21	11	76	37	12	73	24	146	50	29	4	83	182	20	8	210
S19	S14	22	13	29	45	109	55	21	43	51	170	52	5	5	62	33	26	8	67
	S15	24	42	75	42	183	54	62	34	22	172	27	46	4	77	81	22	14	117
	S16	14	21	30	7	72	71	24	98	17	210	37	6	4	47	73	28	7	108
	S17	49	39	19	16	123	22	19	25	17	83	26	7	4	37	57	10	11	78

# ANALYSIS

- A detailed analysis process follows the study
- All the information should be represented in a proper way (mainly, spreadsheets)

Version	Day	Hour	Age	Gender	Experience FM	Yrs Alloy	Most Attractive	Most Relevant	Which would use	Recommend
S01			Anima		Alloy		Anima		Alloy	
S02			Task 01		Task 01		Task 02		Task 02	
S03										
S04										
S05										
S06	S01									
S07	S02									
S08	S03	S01	0	0	0	1	0	0	0	1
S09	S04	S02	0	0	0	0	0	1	1	0
S10	S05	S03	0	0	0	0	0	0	0	1
S11	S06	S04	0	0	0	0	0	0	0	0
S12	S07	S05	0	0	0	0	0	0	0	1
S13	S08	S06	0	0	0	1	0	0	0	1
S14	S09	S07	0	0	0	0	0	0	0	1
S15	S10	S08	0	0	0	0	0	0	0	1
S16	S11	S09	0	0	0	0	0	0	0	1
S17	S12	S10	0	0	1	0	0	0	0	1
S18	S13	S11	0	0	0	1	0	0	0	1
S19	S14	S12	0	0	0	0	1	0	0	1
	S15	S13	0	0	0	0	0	0	0	0
	S16	S14	0	0	0	1	0	1	0	0
	S17	S15	0	0	0	1	0	0	0	0
		S16	0	0	0	0	0	0	0	0
		S17	0	0	0	0	0	0	0	0

# ANALYSIS

- A detailed analysis process follows the study
- All the information should be represented in a proper way (mainly, spreadsheets)

	Version	Day	Hour	Age	Gender	Experience FM	Yrs Alloy	Most Attractive	Most Relevant	Which would use	Recommend
S01											
S02											
S03											
S04											
S05											
S06											
S07											
S08											
S09											
S10											
S11											
S12											
S13											
S14											
S15											
S16											
S17											
S18											
S19											
S20											
S21											
S22											
S23											
S24											
S25											
S26											
S27											
S28											
S29											
S30											
S31											
S32											
S33											
S34											
S35											
S36											
S37											
S38											
S39											
S40											
S41											
S42											
S43											
S44											
S45											
S46											
S47											
S48											
S49											
S50											
S51											
S52											
S53											
S54											
S55											
S56											
S57											
S58											
S59											
S60											
S61											
S62											
S63											
S64											
S65											
S66											
S67											
S68											
S69											
S70											
S71											
S72											
S73											
S74											
S75											
S76											
S77											
S78											
S79											
S80											
S81											
S82											
S83											
S84											
S85											
S86											
S87											
S88											
S89											
S90											
S91											
S92											
S93											
S94											
S95											
S96											
S97											
S98											
S99											
S100											

Scenario 1												Scenario 2								
Anima						Alloy						Anima								
Subject_ID	First_Name	Age	Gender	MD_Rating	PD_Rating	TD_Rating	Performance_R	Effort_Rating	Frustration_Rating	MD_W										
s01	s01	21	Male	0	0	15	70	35	20	0.3333										
s02	1	21	Female	30	0	10	15	20	0	0.2666										
s03	1	22	Male	11	0	15	3	22	0	0.2666										
s04	1	22	Male	40	40	22	6	30	20	0.3333										
s05	1	21	Male	15	0	0	20	52	0	0.2666										
s06	1	21	Female	40	0	0	10	11	5	0.2666										
s07	1	27	Male	20	20	30	20	25	10	0.1333										
s08	1	22	Female	30	3	6	90	30	0	0.2										
s09	1	22	Female	30	40	11	21	26	0	0.2										
s10	1	22	Male	30	5	11	10	25	6	0.1333										
s11	1	21	Male	11	15	70	11	14	11	0.2										
s12	1	21	Male	40	30	40	30	50	35	0.3333										
s13	1	21	Female	65	5	50	74	65	66	0.3333										
s14	1	22	Female	27	10	30	19	40	0	0.2										
s15	1	21	Male	20	0	20	15	50	15	0.3333										
s16	1	21	Male	32	20	23	13	15	5	0.2										
s17	1	21	Female	34	11	50	29	36	22	0.2666										
s18	1		Male	36	27	29	24	45	32	0.3333										
s19	1		Male	45	40	74	10	59	11	0.2										
s01	2		Male	0	30	15	60	35	30	0.3333										
s02	2	1	Female	55	44	20	22	50	55	0.2666										



# ANALYSIS

- A detailed analysis process follows the study
- All the information should be represented in a proper way (mainly, spreadsheets)

S01

S02

S03

S04

S05

S06

S07

S08

S09

S10

S11

S12

S13

S14

S15

S16

S17

S18

S19

Version

Day

Hour

Age

Gender

Experience FM

Yrs Alloy

Most Attractive

Most Relevant

Which would use

Recommend

Anima

Alloy

Anima

Alloy

Task 01

Task 01

Task 02

Task 02

Scenario 1

Scenario 2

Anima

Alloy

Anima

Subject\_ID

First\_Name

Age

Gender

MD\_Rating

PD\_Rating

TD\_Rating

Performance\_R

Effort\_Rating

Frustration\_Rating

MD\_W

s01

s01

21

Male

0

0

15

70

35

20

0.3333

s02

0.2666

s03

0.2666

s04

0.3333

s05

0.2666

s06

0.2666

s07

0.1333

s08

0.2

s09

0.2

s10

0.1333

s11

0.2

s12

0.3333

s13

0.3333

s14

0.2

s15

0.3333

s16

0.2

s17

0.2666

s18

0.3333

s19

0.2

s01

0.3333

s02

0.2666

# ANALYSIS

---

- The data should then be processed in order to answer the fundamental question:
  - Is there a difference in performance, while using both tools?
- Depending on the test, different approaches are used
  - Time & errors: mean times are directly compared
  - NASA TLX & UEQ: Use provided tools for calculating the results

# ANALYSIS

---

- But are the differences significative?
  - Even if a tool performs 50% better than the other one, for a given dimension, it doesn't mean that there is a significative difference
- Depending on the size and distribution of the sample, the T-Test, and the Wilcoxon Signed-Rank Test can be used
- Kolmogorov-Smirnov is used to check for data normality

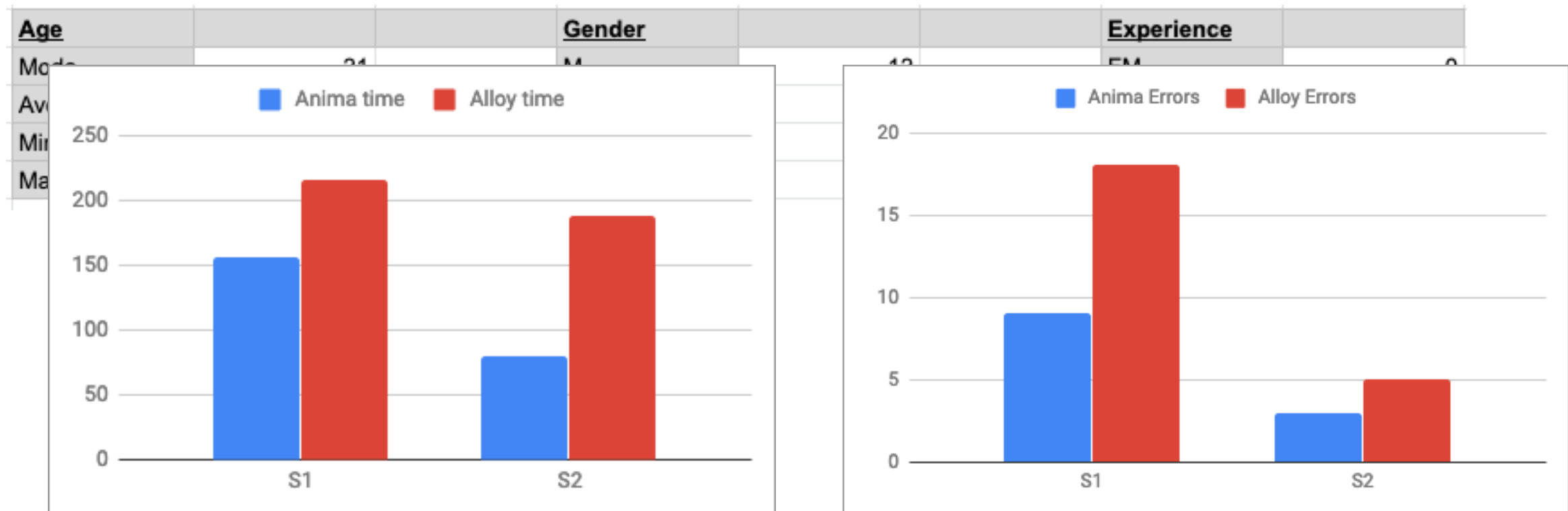
# ANALYSIS

.....

- Finally, the data is represented in a meaningful representation, such as charts

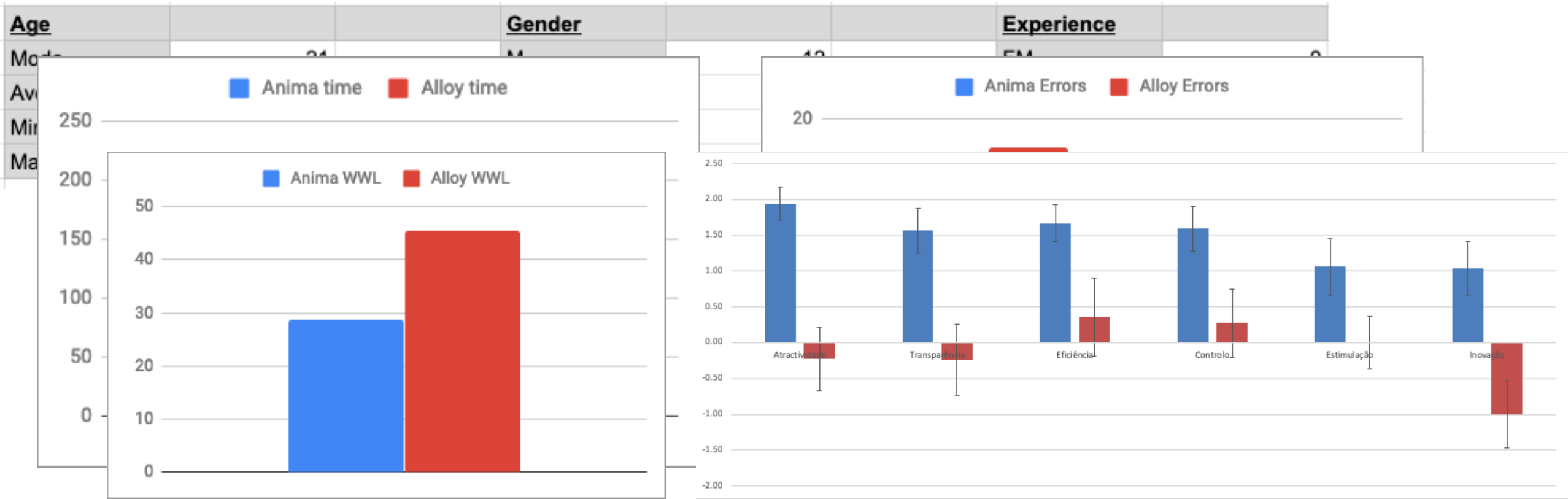
<u>Age</u>			<u>Gender</u>			<u>Experience</u>	
Mode	21		M	12		FM	0
Average	22,29665072		F	7		Alloy	0
Min	21						
Max	33						

- Finally, the data is represented in a meaningful representation, such as charts



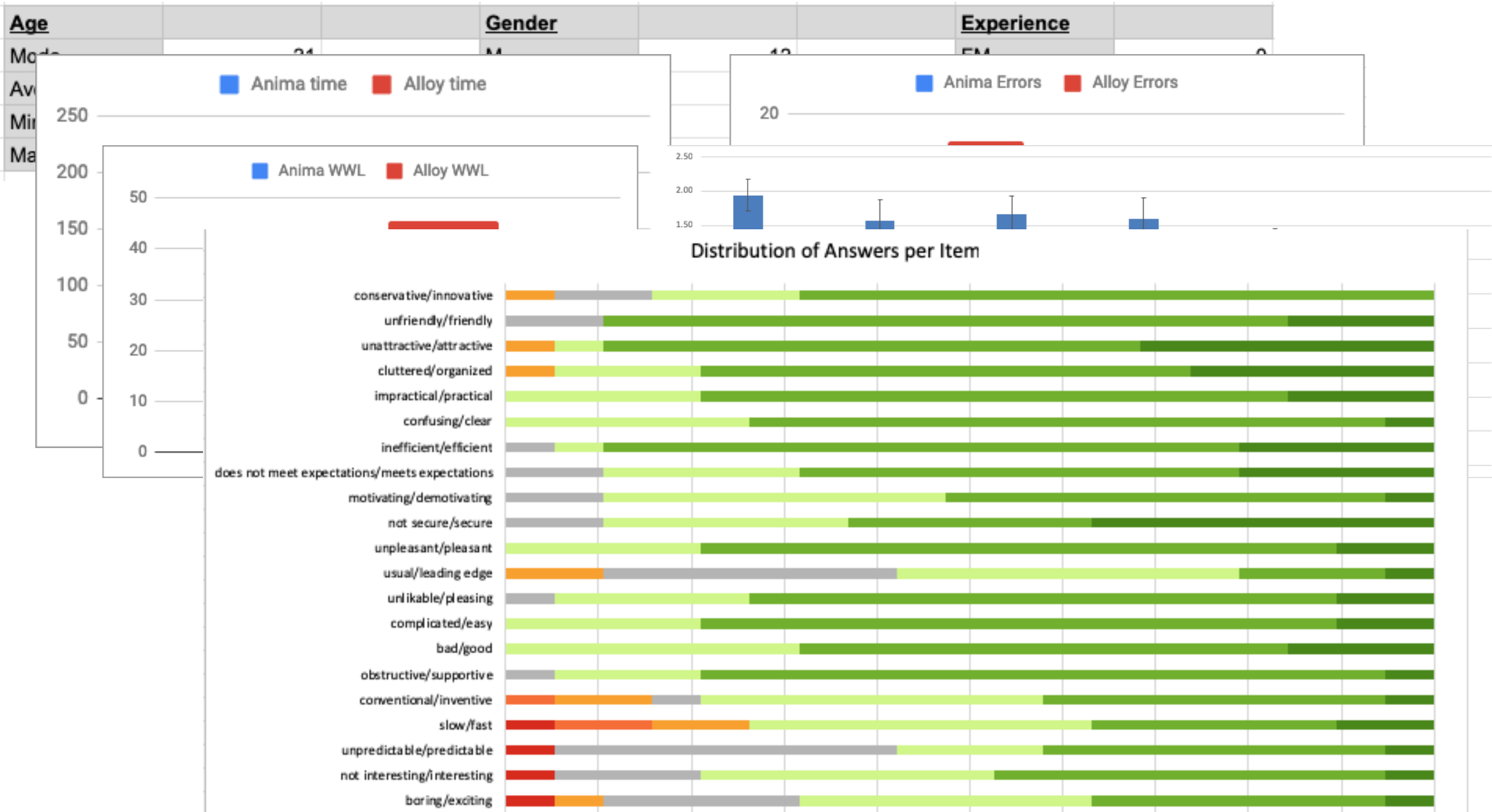
# ANALYSIS

- Finally, the data is represented in a meaningful representation, such as charts



# ANALYSIS

- Finally, the data is represented in a meaningful representation, such as charts





# FINAL RESULTS

---

- Regarding our study, the results have shown that:
  - The users were faster when using Anima
  - The users performed less errors with Anima
  - In some cases, some users did not consider Alloy representations worse
  - Animations and consistency between states were pointed by the participants as the biggest improvements
    - The same is true for the images.



# EMPIRICAL USER STUDIES

---

*Evaluation of the Anima tool*  
*Rui Couto • José C. Campos*