



**HAL**  
open science

# Making sense of learner behavioral, cognitive and demographic characteristics to improve learner modeling

Benoît Choffin, Alice Latimier, Niluphar Ahmadi

## ► To cite this version:

Benoît Choffin, Alice Latimier, Niluphar Ahmadi. Making sense of learner behavioral, cognitive and demographic characteristics to improve learner modeling. International Congress on Technologies in Education, May 2019, Paris, France. , 2019. hal-02139264

**HAL Id: hal-02139264**

**<https://hal.science/hal-02139264>**

Submitted on 24 May 2019

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License



# Making sense of learner behavioral, cognitive and demographic characteristics to improve learner modeling

**Benoît Choffin**  
LRI/CentraleSupélec  
Univ. of Paris-Saclay, Gif-sur-Yvette  
benoit.choffin@lri.fr

**Alice Latimier**  
Département d'Etudes Cognitives (LSCP)  
Ecole Normale Supérieure, Paris  
alice.latimier@ens.fr

**Niluphar Ahmadi**  
LP3C, Synlab  
University of Rennes 2, Rennes  
niluphar.synlab@gmail.com

## Context

- Learner modeling techniques are crucial for providing a personalized and efficient adaptive instruction to learners [3]
- Learning management systems (LMS) automatically record online log-file data : e.g. number of clicks or minutes learners spent on a certain task
- Research in the field of educational data mining used log-files to identify learning strategies and classify learners with respect to their strategy use

→ **Log data are objective information on the use of learning strategies [1].**

For example, Fang et al., 2018 [2] used log-file data to improve adaptivity in CSAL AutoTutor thanks to a better characterization of the students' learning behaviors. They used cluster analysis (k-means + HCA) to create clusters of learners based on interaction logs (**253** learners) from CSAL AutoTutor.

→ **They distinguished 4 clusters of learners** : “**proficient readers**”, “**struggling readers**”, “**conscientious readers**”, “**disengaged readers**”.

## The present study

### Objectives

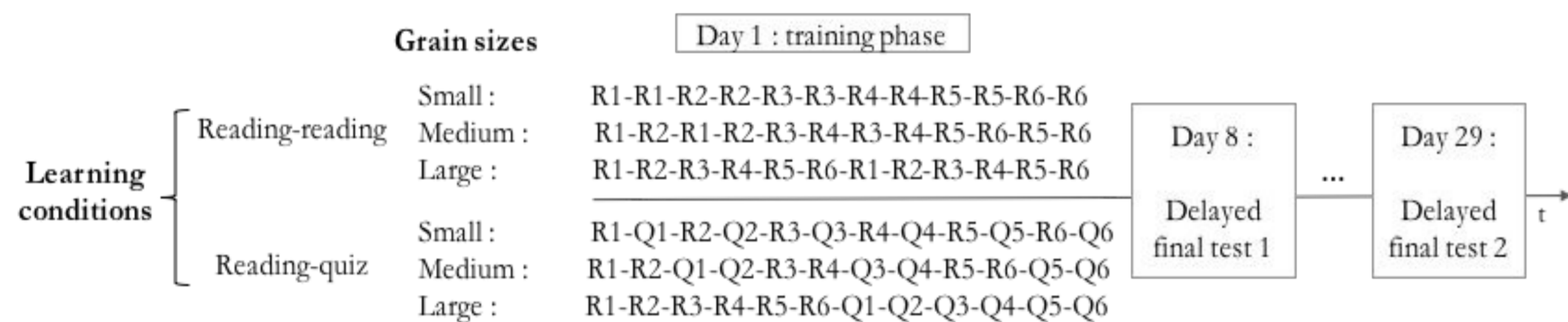
- 1) **Shedding light on underlying behavioral learner patterns**
- 2) **Improving learner predictive modeling to better tailor adaptive online tutoring systems**

### Main research questions

- How do learners interact with a digital learning platform when they are assigned to specific learning strategies?
- Do learning outcomes depend on learners' interactions with the platform (training performances, times on contents...)?
- Do learning outcomes depend on learners' individual features (socio-demographic data)? And to what extent?
- Are the learning strategies effects/results driven by these predictors, variables? What are the interactions?

## Methods

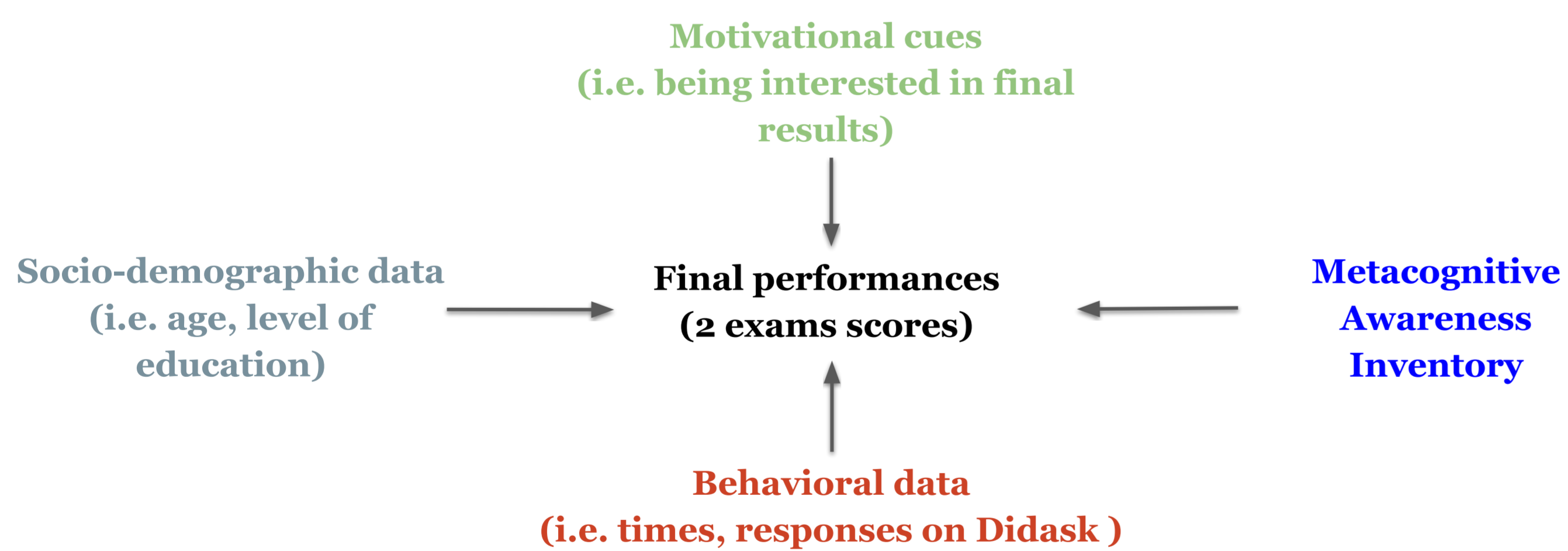
- We used log data from a **learning experiment** conducted with the digital platform Didask. It aimed at comparing **different grain sizes** of learning contents (small/medium/large) and **reviewing strategies** (retrieval practice/reading) on students' performance at two delayed tests.



**DESIGN:** 2 between-subjects variables and 2 within-subjects variables:

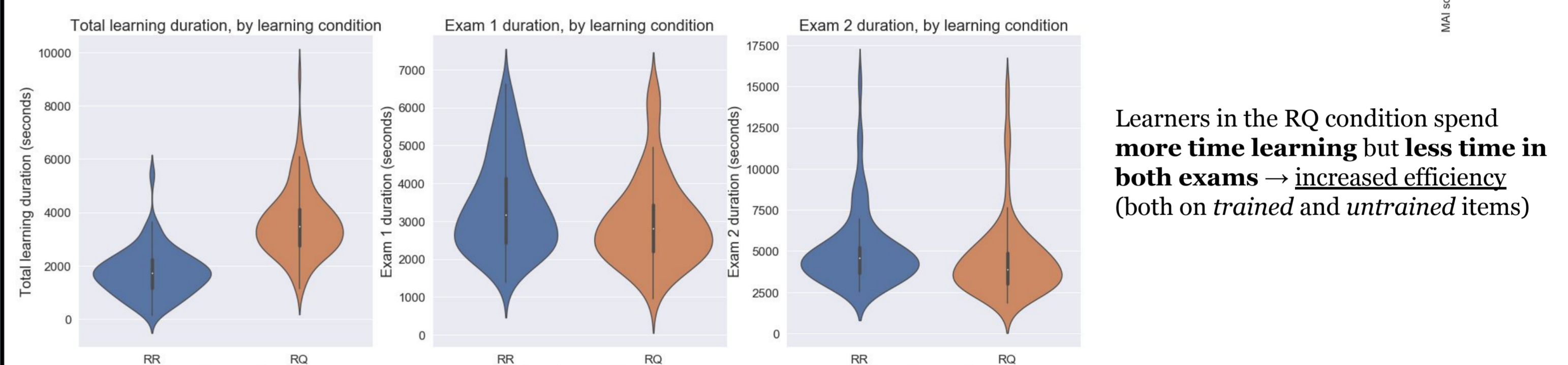
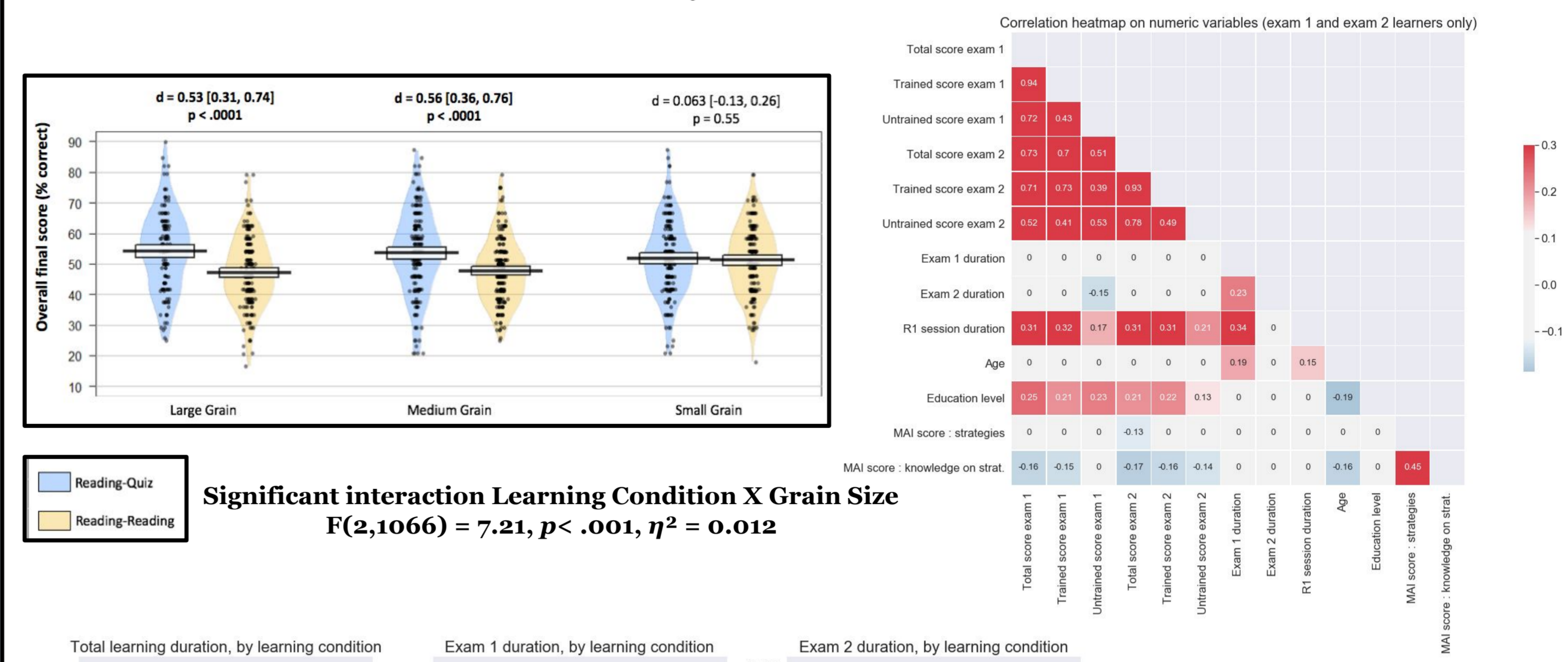
**2 Learning Conditions X 3 Grain Sizes X 2 Question Types X 2 Retention Intervals**

- After aggregation, four types of variables were considered in our analysis:



- Data preprocessing:
  - Data cleansing and variables aggregation (e.g. sum of durations, mean of quiz scores);
  - Specific attention to *duration data*: presence of outliers → replaced by *mean* (over chapter unit if relevant) when above a threshold (determined after histogram screening);
- After preprocessing: **251 students** who passed both exams and **294 students** who passed only exam 1
  - Balanced across learning conditions and grain sizes: RR (52%) vs RQ (48%), SGS (34%) vs MGS (35%) vs LGS (31%)
- Metacognitive Awareness Inventory (MAI) → **validated** in our population
  - 10 items were kept in total
  - Factor analysis → 2 *dimensions* (strategies and general knowledge on metacognition)

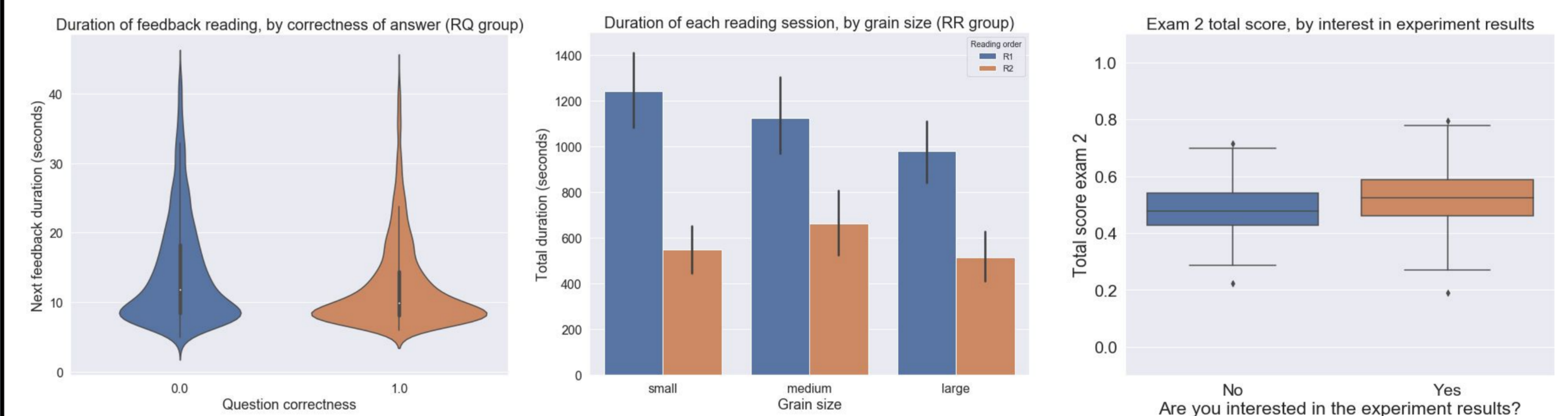
## Data analysis: main results



When RQ learners answer an item *correctly* during the learning phase, they tend to spend **less time reading the feedback**, even though it can be beneficial to them.  
→ after incorrect: **14,43 seconds** on avg.  
→ after correct: **12,34 seconds** on avg.

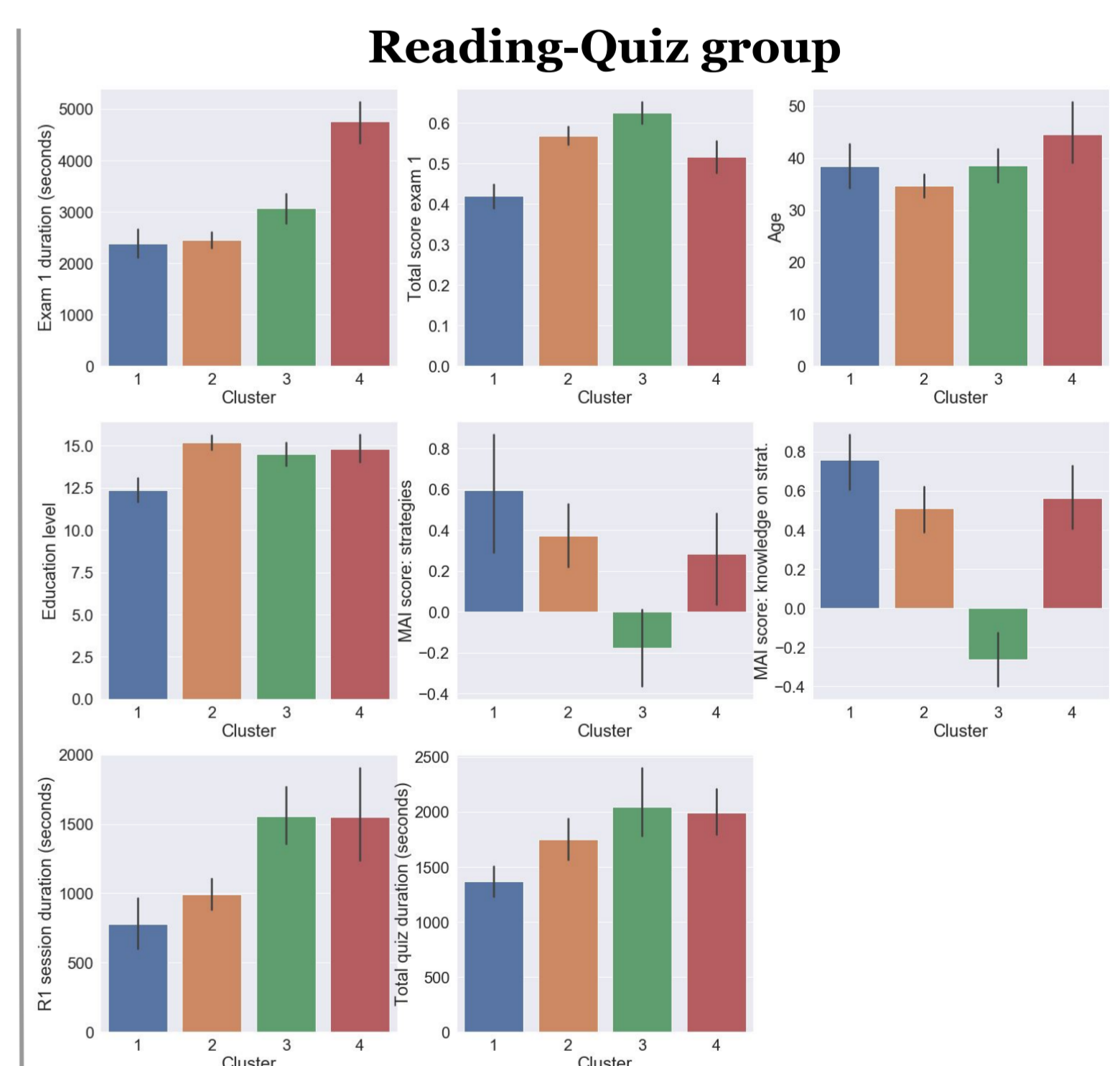
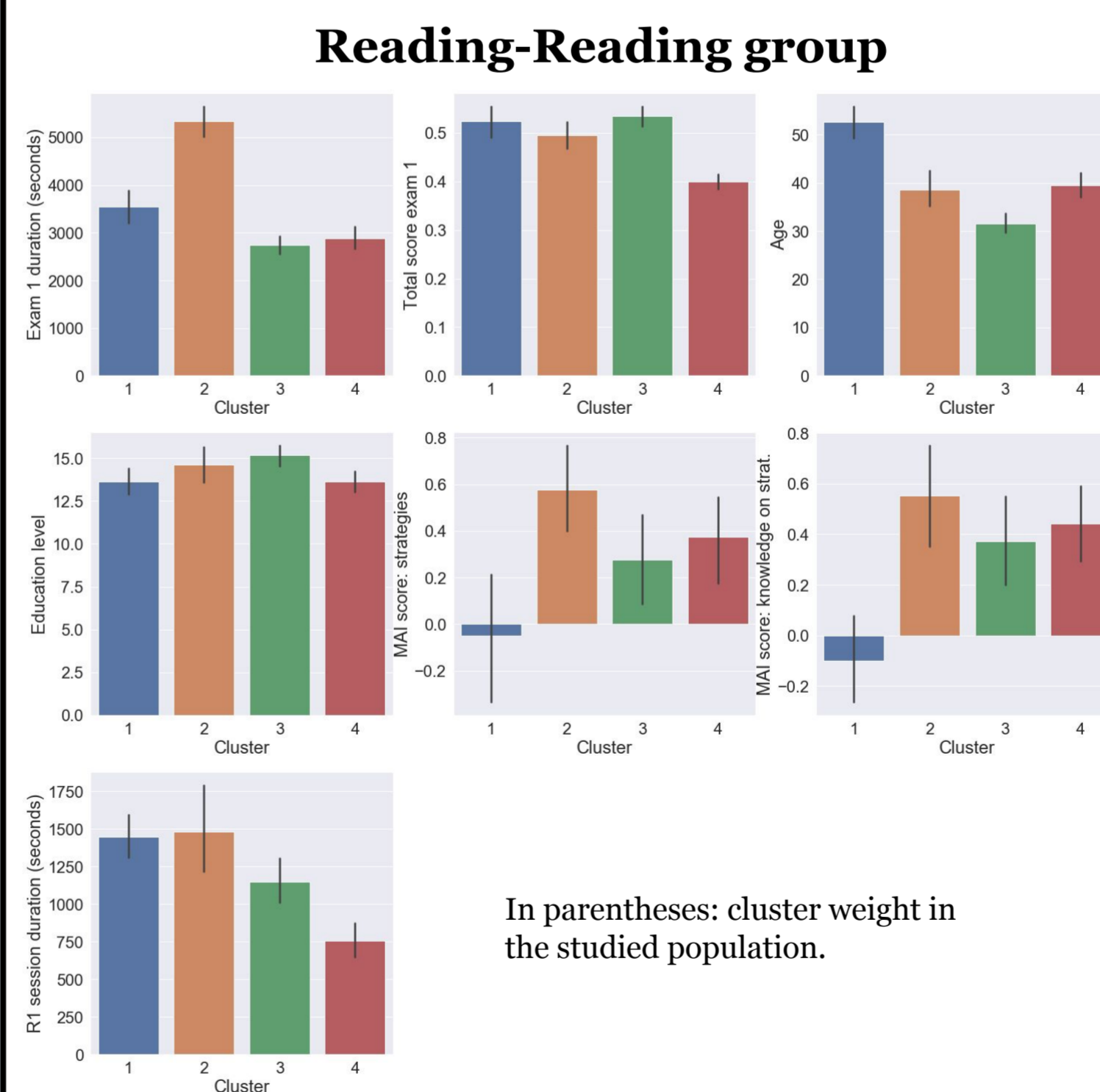
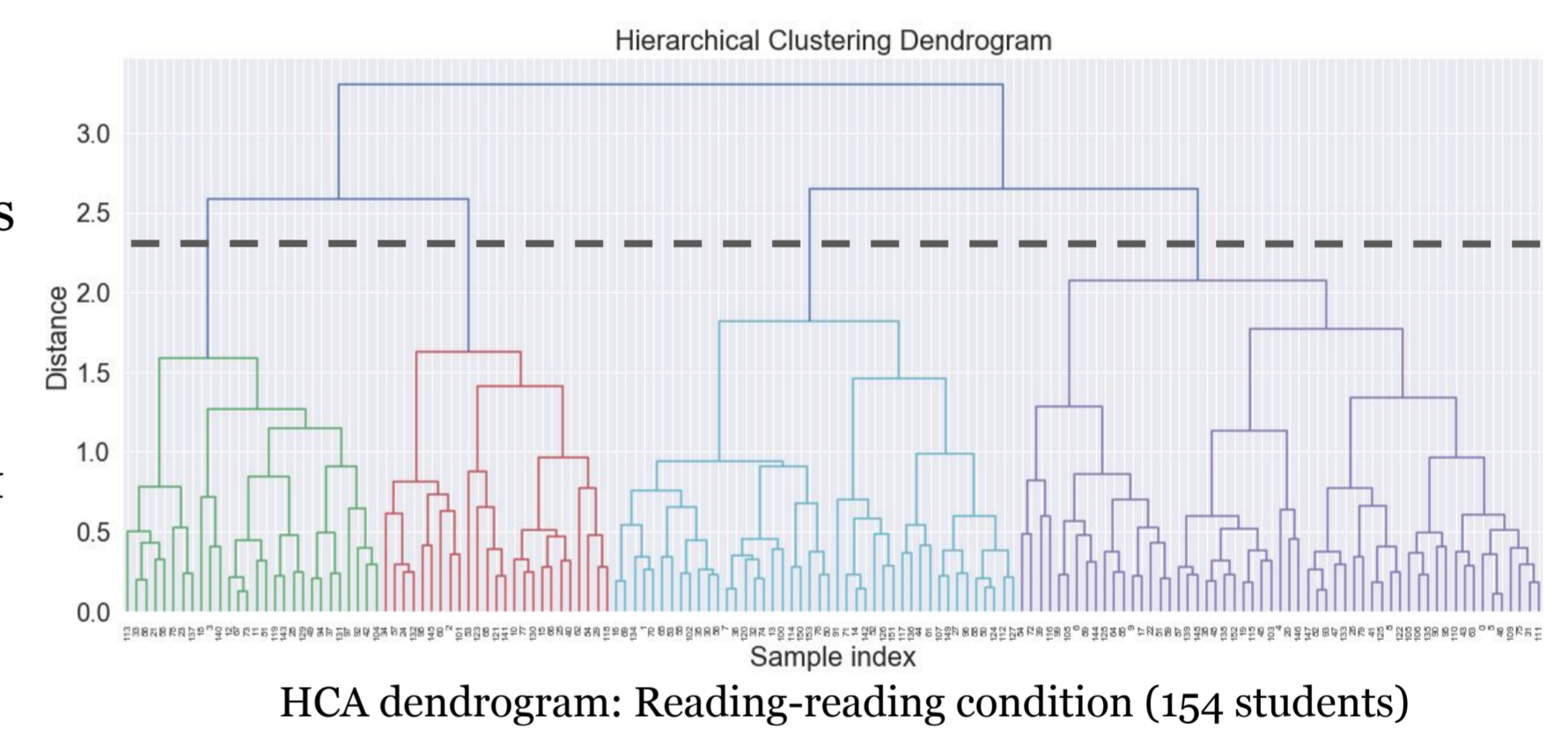
**Illusion of mastery:** RR learners spend less time on R2 than on R1.  
**Progressive disengagement:** as the grain size increases, time spent on R1 decreases (possibly due to motivation loss).

**Motivational cue:** students interested in the results of the experiment have higher scores on both exams than uninterested students (*intrinsic motivation*).



## Clustering: method and results

- **Method: Hierarchical Cluster Analysis (HCA)** → partitions the population of students into homogeneous groups (based on their characteristics)
- Learners in a given cluster have **similar features and behavioral patterns**
- Performed on the 251 students (exam 1 & 2), separately on each learning condition (RR/RQ) → allows to compare clusters between conditions



- 4 clusters:
1. **The skilled learners** (18%)
  2. **The conscientious learners** (16%)
  3. **The efficient learners** (29%)
  4. **The disengaged learners** (37%)

- 4 clusters:
1. **The illusioned learners** (21%)
  2. **The efficient and confident learners** (38%)
  3. **The conscientious but under-confident learners** (25%)
  4. **The struggling learners** (16%)

## Highlights & Perspectives

- Possible to use data from an e-learning platform to do post hoc analyses (could be interesting to use log file from a real learning context) → good support for **formulating new hypotheses** and **guiding future experiments**
- The IV (learning conditions and grain size) are not the only variables that explained the results, **effects are modulated by uncontrolled variables**

- Integrating both numeric and categorical variables inside the cluster analysis (e.g. *AFDM*)
- Using a *mediation model* to look at direct and indirect effects of the covariates on the exam grades
- More complex modelling: *Bayesian network* (uncover more complex dependencies)

### References

1. Theobald, M., Bellhauser, H., & Imhof, M. (2018). Identifying individual differences using log-file analysis: Distributed learning as mediator between conscientiousness and exam grades. *Learning and Individual Differences*, 65, 112-122.
2. Fang, Y., Shubeck, K., Lippert, A., Cheng, Q., Shi, G., Feng, S., ... & Frijters, J. (2018). Clustering the Learning Patterns of Adults with Low Literacy Skills Interacting with an Intelligent Tutoring System. *International Educational Data Mining Society*.
3. Baker, R. (2016). Using learning analytics in personalized learning. *Handbook on personalized learning for states, districts, and schools*, 165-174.