



Fachhochschule Nordwestschweiz
Hochschule für Angewandte Psychologie

Sequencing Immersive Virtual Reality and Physical Laboratory Work for Titration Training: Effects on Learning Outcomes

MASTER-ARBEIT

2025

Autorin / Autor
Chalupny, Urs

Begleitperson
Christ, Oliver

Praxispartner*in
FHNW
Kontaktperson Christ, Oliver

Sequencing Immersive Virtual Reality and Physical Laboratory Work for Titration Training: Effects on Learning Outcomes

Urs Chalupny^{1,*}

¹Angewandte Psychologie, Fachhochschule Nordwestschweiz (FHNW), Olten, Switzerland

Correspondence*:

Urs Chalupny
urs.chalupny@students.fhnw.ch

2 ABSTRACT

3 Higher education institutions preparing students for careers in the chemical industry face mounting
4 challenges in providing quality laboratory training due to increasing student enrollment, limited resources,
5 and safety concerns, prompting exploration of immersive virtual reality as a potential solution for
6 chemistry education. This study investigated the effectiveness of different sequences for integrating
7 virtual reality simulation with traditional laboratory work in teaching titration skills to undergraduate
8 Life Sciences students. Ninety-one students participated in a mandatory course where they completed
9 titration procedures in both virtual reality and physical laboratory environments in counterbalanced order,
10 with knowledge assessed through tests at three timepoints and psychological factors measured using
11 the Cognitive Affective Model of Immersive Learning framework. Results revealed differences in learning
12 outcomes between virtual reality and physical laboratory conditions, nor did the sequence of exposure
13 (virtual reality first versus laboratory first) affect knowledge acquisition. Students demonstrated high
14 baseline knowledge (75% correct) with modest but equivalent gains regardless of learning modality. Path
15 analysis confirmed only one significant relationship from the theoretical model: higher sense of presence
16 in virtual reality was associated with reduced cognitive load. Individual differences among participants
17 accounted for 76% of the variance in learning outcomes, while experimental conditions explained less than
18 1%. These findings demonstrate that immersive virtual reality can serve as a pedagogically equivalent
19 alternative to traditional laboratory training for procedural chemistry skills, supporting its potential as
20 a scalable complement to physical laboratory experiences when resources are constrained or access is
21 limited.

22 **Keywords:** Virtual Reality, higher education, Laboratory, Immersive Learning

1 INTRODUCTION

23 The development and mastery of practical laboratory skills remain a cornerstone of education in
24 chemistry and related disciplines, crucial for equipping future professionals for the demands of the
25 labor market. A growing emphasis is placed on shifting from traditional, content-focused curricula
26 towards competency-based education (Campbell et al., 2022; European Commission and PwC, 2020).
27 Laboratory environments are pivotal in this shift, fostering not only declarative and procedural
28 knowledge but also transversal skills and higher-order thinking, enabling a holistic educational

29 approach (Agustian et al., 2022). The recent COVID-19 pandemic underscored the value of hands-on
30 experience, as limited access to physical labs negatively impacted perceived learning outcomes and
31 confidence in practical skills (Finne et al., 2022; Hyde, 2025). However, higher education institutions
32 are facing major challenges due to steadily increasing student numbers. In Switzerland alone, the
33 number of bachelor's students at universities of applied sciences rose by approximately 30% between
34 2011 and 2021 (Wolter et al., 2023), a trend that is also observed globally (UNESCO, 2024). This
35 development poses significant challenges for institutions aiming to maintain or even expand this
36 form of hands-on training. In the case of laboratory education, it also affects access to lab facilities.

37 In this situation, providing high-quality learning opportunities in traditional laboratories faces
38 significant hurdles and discussions even among specialists (Bretz, 2019). Restricted lab space, large
39 class sizes (Tee et al., 2018), limited availability of qualified personnel, the complexity of effective
40 pedagogical implementation (Agustian et al., 2022), substantial material costs, and inherent safety
41 concerns (Bennett and O'Neale, 1998) collectively challenge educators. Since the early days of
42 education, when direct experiences essential for learning (Meyer-Drawe, 2003) are not possible,
43 instructional methods have been employed to convey declarative or procedural knowledge or to
44 foster the development of competences (Zierer and Seel, 2012). In this context, virtual simulations
45 have increasingly come into use in recent years. Virtual Reality (VR) represents one specific form
46 along the reality–virtuality continuum (Milgram et al., 1995), in which reality is largely replaced
47 by a simulation. In VR, this replacement is primarily achieved through the use of head-mounted
48 displays (HMD), which block the user's visual access to the physical environment and replace
49 it with a stereoscopic two-dimensional image. In our case, we attempted to create a laboratory
50 simulation with the aim of replicating the physical environment, allowing learners to engage in
51 the intended experiences within a virtual environment. These learning experiences can be accessed
52 independently of the physical location of the educational institution and can be made available
53 asynchronously—essentially anywhere that the necessary technological infrastructure is available.
54 The previously described challenges associated with limited access to laboratory training—along
55 with the core educational intention—could therefore be largely addressed through the use of VR.

56 Immersive virtual reality (iVR)—digitally created environments offering interactive and “real-
57 feeling” experiences (Bailenson, 2018)—presents a promising technological avenue to address several
58 of these challenges. Specifically, iVR can mitigate ongoing material costs, requires less physical
59 space, allows virtually unlimited repetition of experiments (Dunnagan et al., 2020), and significantly
60 reduces safety risks (Bøg Petersen et al., 2022; Concannon et al., 2019) and can give access to
61 learning experiences otherwise too dangerous or expensive (Tee et al., 2018). Some studies show non-
62 inferior learning outcomes compared to traditional labs. This makes iVR a potentially cost-effective
63 equivalent (Dunnagan et al., 2020). Successful adoption requires careful consideration of virtual
64 environment design, effective instructional pathways, and the development of necessary competencies
65 among both educators and students (Radianti et al., 2020). Even though the design elements required
66 to create effective learning experiences—both in iVR and in physical laboratories—are relatively
67 well researched, the question remains as to how these two environments can be optimally aligned
68 to ensure that learning and the transfer of skills from one setting to the other actually take place,
69 rather than being hindered.

70 Titration, a fundamental quantitative chemical analysis technique used to determine the
71 concentration of an analyte, serves as a prime example of a core laboratory competency (Sheppard,
72 2006). Mastering titration demands a blend of skills: students must understand the underlying

73 chemical principles (declarative knowledge), correctly execute the procedural steps using specific
74 equipment (procedural knowledge), and apply mathematical skills for accurate calculation. This aligns
75 with [Johnstone \(1993\)](#) framework, requiring integration of macroscopic observations, submicroscopic
76 understanding, and representational competence (symbols, equations) and makes the titration an
77 excellent object for learning different skills and comparing their acquisition in different environments.
78 The specific demands differ between environments: physical labs emphasize fine motor skills and
79 safety diligence, while virtual labs require initial proficiency with the iVR interface before engaging
80 with the chemical procedure itself ([Johnson-Glenberg et al., 2023](#)). These subtle differences between
81 learning environments are highly relevant when considering the transfer of the competences to be
82 developed. A key question is whether the similarities in the actions performed are sufficient to enable
83 the transfer of procedural titration knowledge across environments. Research on transfer emphasizes
84 underlying principles, identical elements, variability in stimuli, and altered learning conditions as
85 critical factors ([Rahman, 2020](#)). In this regard, learning through different environments—but with a
86 shared core task—may prove to be particularly effective.

87 The potential of iVR in education has been explored for decades, aiming to enhance motivation,
88 engagement, and learning ([di Lanzo et al., 2020](#)). However, research reveals a heterogeneous
89 landscape regarding its effectiveness ([Won et al., 2023](#)). While iVR often boosts affective outcomes
90 like motivation, engagement, interest, and presence ([Liu et al., 2023](#); [Makransky et al., 2021](#); [Parong
91 and Mayer, 2018](#)), its impact on cognitive learning outcomes is ambiguous. Several meta-analyses and
92 systematic reviews suggest that iVR can be particularly advantageous for learning procedural skills
93 compared to less immersive methods like presentations or lectures ([Conrad et al., 2024](#); [Hamilton
94 et al., 2021](#)). Studies focusing on procedural tasks often report positive effects ([Concannon et al.,
95 2019](#); [Hamilton et al., 2021](#)), and VR training has shown potential as a precursor to real-world task
96 performance, indicating successful transfer in different contexts ([Hamilton et al., 2021](#)) and especially
97 in the laboratory ([Bøg Petersen et al., 2022](#)). The fidelity (physical and cognitive similarity to the
98 real task) of the simulation is considered crucial for facilitating this transfer between the simulation
99 and the real-world task ([Bøg Petersen et al., 2022](#); [Johnson-Glenberg et al., 2023](#); [Levac et al.,
100 2019](#)). Although a trend appears to be emerging for developing Skills in VR, these findings must be
101 interpreted with caution, as the studies differ considerably in their methodologies, the interventions
102 used, and the operationalization of learning outcomes ([Liu et al., 2023](#); [Radianti et al., 2020](#)).

103 Conversely, when acquiring declarative knowledge, iVR has sometimes led to poorer outcomes
104 compared to traditional methods like slideshows or text ([Bøg Petersen et al., 2022](#); [Makransky and
105 Petersen, 2019](#); [Parong and Mayer, 2018](#)). A review by [Tusher et al. \(2024\)](#) stated as well, that the
106 influence of VR on cognitive and affective learning outcomes appeared to be “subdued” compared to
107 its beneficial effects on motor skills. Numerous empirical investigations ([Makransky and Petersen,
108 2019](#); [Parong and Mayer, 2021a,b](#)) have drawn upon Sweller’s Cognitive Load Theory ([Sweller et al.,
109 2011](#)) to explain how environmental factors influence learning outcomes, particularly in the acquisition
110 of declarative knowledge. The richness and interactivity of iVR environments, while potentially
111 engaging, might introduce extraneous cognitive load (processing irrelevant information) or increase
112 intrinsic load (complexity of interaction), overwhelming the learner’s working memory capacity and
113 hindering knowledge construction ([Bøg Petersen et al., 2022](#); [Makransky and Petersen, 2019](#); [Parong
114 and Mayer, 2021a](#)). This raises the question of optimal sequencing of the learning experiences; might
115 prior real-world experience reduce the cognitive load in VR or vice-versa? Furthermore, the novelty
116 effect, where unfamiliarity with the technology itself consumes cognitive resources, can negatively
117 impact initial learning sessions ([Hamilton et al., 2021](#); [Parong and Mayer, 2021a](#)). Therefore, effective

118 instructional design within VR, applying principles like signaling (highlighting key information)
119 and spatial/temporal contiguity (presenting related information together) derived from multimedia
120 learning theory (Peeters et al., 2023), is critical but often under-implemented (Liu et al., 2023;
121 Radianti et al., 2020).

122 To better understand learning mechanisms within immersive environments like iVR, the Cognitive
123 Affective Model of Immersive Learning (CAMIL), developed by Makransky and Petersen (2021),
124 provides a comprehensive framework for understanding the interplay of technological, cognitive, and
125 affective factors that influence learning outcomes. It englobes technological factors, like immersion,
126 that facilitate the feeling of presence in the virtual environment and enable to some degree the
127 agency of the user and the different possibilities to interact and to change and control objects in
128 the simulation. These two concepts, presence and agency, are defined as affordances of iVR and
129 are the factors that influence the affective components of the model. Previous studies found that
130 a high sense of presence was related with higher levels of situational motivation (Makransky and
131 Lilleholt, 2018). Another important cognitive factor is the extraneous cognitive load (CL) depending
132 on the virtual environment and the degree of stimuli and cues that have to be treated by the user
133 when navigating and interacting with the simulation. This CL, as theorized in the CLT, would
134 occupy cognitive resources that would be otherwise free and accessible for the learning tasks. So
135 the design of the virtual environment should be optimized in regards to minimize distractions and
136 maximize cues to engage in the learning task. The CAMIL Model states that a higher motivation
137 and lower CL should have a directed influence on learning (Petersen et al., 2022). The model also
138 states that external factors like low usability scores, measured eventually with the System Usability
139 Scale (Lewis, 2018) or personality-related factors such as technology affinity and behavioral intention
140 toward technology use, as measured by the TAEG questionnaire (Karrer-Gauß et al., 2009) could
141 impair learning with iVR (Makransky and Petersen, 2021).

142 Given that the CAMIL Model describes external factors as important determinants of the learning
143 process in iVR, this study seeks to elucidate how these influences manifest and interact in educational
144 contexts. These different relationships between technological and intrapersonal factors are crucial
145 when trying to understand the emergence of learning in VR and one aim of this study is to better
146 understand these possible influences. We assume that the acceptance of the use of technology in the
147 specific case could be a predictor of presence in iVR. Links between acceptance and performance
148 when using VR devices were shown (Barrett et al., 2023; Guo et al., 2025) and we hypothesize that
149 the perception of the technology should have a direct influence on the presence in iVR.

150 **H1:** Technology affinity (TA) should have a positive relation to perceived presence in iVR. TA is
151 defined as a personality trait that manifests in a positive attitude, enthusiasm, and confidence toward
152 technology (Karrer-Gauß et al., 2009) and can be measured with a short questionnaire developed by
153 Karrer-Gauß et al. (2024).

154 In its original formulation, the CAMIL model posited that presence influences extraneous cognitive
155 load, specifically suggesting that presence would influence extraneous cognitive load stemming from
156 the environment (Makransky and Petersen, 2021). However, this direct relationship has not yet
157 been empirically confirmed and was questioned in the study by Petersen et al. (2022), who found
158 a negative effect of cognitive load on presence rather than the other way around. This indicates
159 that high cognitive load from interaction decreased users' feelings of presence, challenging the
160 presumed direction of influence in that context. It should be noted, however, that their study
161 involved a learning task focused exclusively on declarative knowledge. Furthermore, while not

empirically demonstrated in their specific VR training scenario, it has been theoretically argued and suggested that factors enhancing subjective experience, such as improved fidelity or relevant cues, can contribute to reducing cognitive load (Cooper et al., 2021). Consistent with this, a systematic review of evidence-based design and pedagogical principles in educational VR environments, drawing on frameworks such as multimedia and generative learning, specifically discusses signaling principles, including the use of textual annotations as an adjunct technique, and references studies linking signaling in VR to cognitive load (Oje et al., 2025). Based on these considerations, and acknowledging the need for further empirical evidence on this specific link, we assume that in the context of our study, a skill-learning task supported by explicit visual and textual cues, the original hypothesis that presence reduces cognitive load may still hold. We posit that a strong sense of presence, facilitated by these cues, could minimize the extraneous cognitive load associated with navigating and performing complex procedures in the virtual environment.

H2: Presence (Pre) should have a negative impact on extraneous CL in iVR.

While empirical findings regarding the direct link between presence and intrinsic motivation are not always directly supported (e.g. Petersen et al., 2022), found physical presence predicted situational interest, which covaried with intrinsic motivation, but not intrinsic motivation directly in their specific model), there is broader support that the use of technologies like VR is associated with a range of positive affective and motivational reactions (Chirico and Gaggioli, 2019). Numerous studies have reported outcomes such as increased motivation (Bohne et al., 2021), self-efficacy (Thisgaard and Makransky, 2017), enjoyment (Cooper et al., 2021), interest (Petersen et al., 2022), and positive affect (Chirico and Gaggioli, 2019) stemming from engagement with VR/AR environments. These positive affective and motivational states are considered important, as theoretical frameworks and empirical evidence suggest they can, in turn, influence learning outcomes.

H3: Presence should have a positive relation to motivation (Mo).

H4: Cognitive Load and motivation should have direct influence on the learning outcome in iVR.

Titration, requiring both procedural execution and conceptual understanding, represents an ideal case study given the mixed findings on iVR effectiveness for different knowledge types. Direct comparisons between VR and real-life training exists (Bøg Petersen et al., 2022; Kaplan et al., 2021), but these studies employed usually a single-session design. Regarding the optimal integration of these modalities within an authentic educational context, the question remains open how to design such learning interventions. Specifically, how does the sequence in which students engage with iVR simulation and traditional laboratory practice impact their learning of a complete procedure like titration? Real laboratory work and learning is crucial (Hyde, 2025), but is there a sequencing or learning condition effect when designing the instruction between the iVR Simulation and the laboratory?

Most studies compare VR against a single alternative or focus on isolated skill components, often outside regular curricula (Hamilton et al., 2021; Radianti et al., 2020). Few have directly compared the sequential application (e.g., VR-first then Lab vs. Lab-first then VR) within an established university course, measuring learning outcomes for the entire procedure. Understanding this sequence effect is crucial for deriving practical guidelines on how best to leverage iVR as a complement to, rather than a replacement for, traditional lab work (Bøg Petersen et al., 2022; Concannon et al., 2019).

204 This study aims to address this gap by investigating the effectiveness of two different sequential
205 learning approaches for teaching titration to undergraduate Life Sciences students within their
206 mandatory “Basic Laboratory Techniques” module. Specifically, we compare the impact on knowledge
207 acquisition when students first perform the titration procedure in an immersive virtual reality
208 simulation followed by the traditional laboratory (VR-Lab sequence), versus performing the
209 traditional laboratory exercise first followed by the VR simulation (Lab-VR sequence).

210 Drawing on cognitive load theory and the principle of scaffolding (Oje et al., 2025), we hypothesize
211 that the VR-Lab sequence may offer advantages over the Lab-VR sequence. In the VR-Lab condition,
212 students can first explore the procedure in a low-stakes environment where errors have no material
213 consequences. This initial exposure may serve as cognitive scaffolding, allowing students to focus
214 on understanding the procedural sequence before managing the additional demands of physical
215 manipulation and safety concerns in the real laboratory. Conversely, the Lab-VR sequence may result
216 in the virtual environment being perceived as redundant after real-world experience, potentially
217 reducing engagement and learning gains.

218 **H5:** Students in the VR-Lab sequence will demonstrate greater knowledge gains compared to those
219 in the Lab-VR sequence, as the virtual environment provides cognitive scaffolding that prepares
220 students for the complex demands of real laboratory work.

221 In summary, this study makes both theoretical and practical contributions to understanding VR-
222 based laboratory education. Theoretically, we extend the CAMIL model by testing its applicability
223 to procedural skill learning in chemistry education, specifically examining whether the relationships
224 between technology affinity, presence, cognitive load, and motivation hold in this applied context.
225 Practically, we address a critical gap in instructional design by investigating how the sequencing
226 of VR and physical laboratory experiences affects learning outcomes—a question of immediate
227 relevance to educators facing resource constraints and seeking evidence-based integration strategies
228 for emerging technologies. These two investigations are complementary: understanding the cognitive-
229 affective mechanisms in VR learning (CAMIL) can help explain why certain sequences might be
230 more effective, while the sequencing comparison provides ecological validity for the CAMIL model
231 in authentic educational settings.

2 METHODS

232 The study was conducted as part of the Bachelor of Science in Life Sciences degree program
233 at the University of Applied Sciences Northwestern Switzerland (FHNW). It was integrated
234 into the foundational module ‘Practical Basic Laboratory Techniques’ (Praktikum Grundlagen
235 Labortechniken), a mandatory course for all students organized during their first semester.

2.1 Participants

237 A total of 92 students initially participated in the study. One participant discontinued the VR
238 condition due to experiencing discomfort and was subsequently removed from the sample. The final
239 sample consisted of 91 participants who completed all stages of the study. Of these, 57 participants
240 (62.6%) identified as female, 33 (36.3%) as male, and one (1.1%) as non-binary. The mean age for the
241 entire sample was 21.82 years. There was no significant difference in mean age between the two largest
242 gender groups ($t_{\text{Welch}}(87.48) = -0.17$, $p = 0.86$, $n_{\text{obs}} = 90$). Regarding the highest educational
243 qualification attained, 36 participants (39.6%) reported a ‘Berufsmatura’ (vocational baccalaureate),

244 44 (48.4%) a ‘Fachmatura’ or ‘gymnasiale Matura’ (specialized or academic baccalaureate), 4 (2.2%)
245 a Federal Diploma of Vocational Education and Training (‘Eidgenössisches Fähigkeitszeugnis’, EFZ),
246 and 2 (2.2%) reported qualifications from higher vocational education. One participant (1.1%) already
247 held a Bachelor’s degree. Two participants (2.2%) selected ‘Other,’ specifying ‘Fachhochschulreife’
248 or ‘Fachmaturität’ (types of university of applied sciences entrance qualifications). Concerning
249 prior experience with VR headsets, 34 participants (37.4%) reported having never used one before.
250 42 participants (46.2%) indicated having very limited prior experience. 13 participants (14.3%)
251 had used a headset multiple times but indicated needing some support, while 2 (2.2%) described
252 themselves as experienced and largely autonomous users. At the time of the study, 3 participants
253 owned a personal VR headset.

254 **2.2 Materials**

255 The welcome, introduction, and all written tests were administered in a lecture hall on the MuttENZ
256 campus, with tables arranged in six rows and eight tables per row, divided by a central aisle. The
257 VR condition was also conducted within this lecture hall. Upon arrival, participants discovered that
258 VR headsets had been strategically placed in pairs on the tables, automatically establishing the
259 two-person teams that would work together for the duration of the experimental session (VR and
260 Lab). Care was taken to maximize the distance between VR stations; consequently, three stations
261 were set up in each half of the room, separated by an unoccupied row of tables. Some research team
262 members were positioned in the front row of the classroom to monitor and support students in
263 the VR condition, observing both the physical classroom environment and the students’ virtual
264 activities via computer displays. Other team members were stationed in the chemistry laboratory to
265 instruct and support students in the traditional lab condition (cf. Figure 1 for an overview of the
266 study procedure).

267 **2.2.1 Virtual environment and materials**

268 The virtual environment replicates a real laboratory and includes all necessary tools and objects
269 for performing the titration of an unknown solution. The simulation was developed in Unity 6,
270 supplemented with objects from the Asset Store and some Models were developed in Blender 3.0.
271 The workspace is situated within a fume hood and features a burette, measuring cylinder, pipette,
272 beakers, funnel, various solutions, a magnetic stirrer, and a pipette bulb. Integrated on the right
273 side of the fume hood was a series of tutorial videos explaining and demonstrating the handling of
274 key objects, grasping and holding techniques, and movement within the virtual space. Adjacent to
275 this, a SharePoint interface provided access to an Excel spreadsheet for recording measured and
276 calculated solution volumes during the experiment. The simulation is structured into three distinct
277 sequential phases. The first phase serves as an orientation period (“arrival”), allowing participants to
278 familiarize themselves with the environment, view optional tutorial videos, and interact freely with
279 the various objects. There were no time constraints or restrictions on interaction during this phase.
280 Participants independently decided when to initiate the second phase and commence the guided
281 titration. During the guided titration phase, instructions were displayed solely as text on the wall of
282 the fume hood. Depending on the current instruction, visual cues (such as yellow outlines, arrows, or
283 flashing indicators) highlighted the corresponding objects or relevant controller buttons. Successful
284 completion of an action was acknowledged by an auditory cue, after which the instruction advanced
285 to the next step. Students were thus guided through the entire titration process, which precisely
286 mirrored the procedural protocol of a titration performed in a real laboratory. Throughout this

287 phase, students were required to transfer three measured burette volumes to the Excel spreadsheet.
288 These entered values were not automatically verified, and no direct feedback on their accuracy was
289 provided within the simulation. Data entry for each volume was confirmed by pressing a designated
290 button, which concluded the specific instruction and displayed the subsequent one. This learning
291 phase concluded with the final completion of the Excel spreadsheet, including the measured volumes
292 and the calculated concentration value. Following the learning phase, a test phase was initiated by
293 pressing another designated virtual button. In this phase, participants had to perform the same
294 sequence of actions again using an unknown solution, transferring the measured values back into
295 the Excel spreadsheet. The experiment concluded with the final transfer of these values. During this
296 phase, no additional guidance or feedback was provided.

297 2.2.2 Laboratory environment and materials

298 The titration procedure was also performed by all students in the university's physical laboratory.
299 Analogous to the virtual lab procedure, students received printed instructions that guided them
300 step-by-step through the experiment. All necessary materials and equipment were provided at the
301 laboratory benches. Students worked in pairs and conducted the experiment under the supervision
302 of instructors. The procedure was carried out in accordance with applicable safety regulations and
303 instructor guidelines. Participants were required to record the measured volumes on a protocol sheet
304 and calculate the concentration of the unknown substance. Instructors were available to answer
305 questions and provide assistance as needed.

306 2.3 Questionnaires and scales

307 All questionnaires were created with the Enterprise Feedback Suite from Tivian and made accessible
308 via QR code. Participants completed the tests on their personal smartphones, which allowed for
309 flexible administration. A standardized questionnaire was employed, structured such that specific
310 sections were unlocked depending on the testing timepoint and condition (VR or Laboratory). Data
311 collection was conducted anonymously, and all participants provided informed consent prior to
312 participation. Participants identified themselves using an anonymous code assigned at the beginning
313 of the study, along with their respective group number.

314 2.3.1 Demographics and prior experience

315 The first assessment included demographic questions regarding age, gender, and highest educational
316 qualification. Participants' level of experience with VR headsets was assessed using the question:
317 'How would you currently rate your experience with using VR headsets?' Response options
318 ranged across five levels, from 'I have never worn a VR headset before' to 'I use VR headsets
319 completely autonomously' (the full German questionnaire is provided in the Appendix). Following
320 the demographic and experience questions, the four items of the German short version of the
321 Technology Affinity Scale (Karrer-Gauß et al., 2024) were presented. This was succeeded by the
322 first administration of the knowledge test on titration.

323 2.3.2 Knowledge Test

324 The knowledge test was developed by the program instructors and aligned with the expected
325 learning objectives following the completion of the titration procedure. To ensure content validity,
326 the test items were selected by three course instructors, checked for alignment with the module's
327 learning objectives, and reviewed for relevance and consistency with the core concepts of titration.

328 **2.3.2.1 Test Content Categories**

329 The test comprised 13 multiple-choice items with multiple correct selections possible per item,
330 assessing different dimensions of titration competency:

331 **Declarative Knowledge** (4 items): Items assessed understanding of fundamental concepts,
332 including:

- 333 • Required chemicals for titration (Item 1)
- 334 • Essential equipment identification (Item 2)
- 335 • Endpoint recognition criteria (Item 6)
- 336 • Theoretical understanding of indicator effects (Item 13)

337 **Procedural Knowledge** (5 items): Items evaluated knowledge of correct laboratory procedures:

- 338 • Correct sequence of titration steps (Item 3)
- 339 • Importance of burette rinsing procedures (Item 8)
- 340 • Necessity of stirring during titration (Item 9)
- 341 • Proper burette filling technique (Item 10)
- 342 • Impact of over-titration on results (Item 11)

343 **Problem-Solving and Application** (4 items): Items required integration of conceptual
344 understanding with practical application:

- 345 • Troubleshooting when burette volume is exhausted (Item 4)
- 346 • Solution strategies for concentration issues (Item 5)
- 347 • Stoichiometric calculations for concentration determination (Items 7 & 12)

348 **2.3.2.2 Scoring System**

349 Items were scored by awarding one point for each correctly selected option and one point for each
350 correctly unselected distractor, resulting in a maximum score of four points per item. Therefore, a
351 maximum total score of 52 points could be achieved. This scoring method was chosen to reduce
352 guessing probability and provide a more nuanced assessment of student knowledge.

353 The item order was randomized to minimize potential order effects. The questions had been used
354 previously in the module and were adapted specifically for this experiment. These adaptations were
355 necessary to standardize the answer format across all questions and to facilitate online administration.
356 Specifically, this required all questions to feature four answer options, and a sorting task (original
357 Item 3) was converted to a multiple-choice format due to limitations of the online testing platform.
358 This adapted version was subsequently reviewed, tested, and ultimately approved by the instructors.

359 **2.3.3 Constructs of the CAMIL Model**

360 Scales for Presence (Pre), Motivation (Mo), as well as intrinsic (CL_int) and extrinsic (CL_ext)
361 cognitive load were drawn from the CAMIL model (Makransky and Petersen, 2021). The items for
362 each scale were described by Petersen et al. (2022) and compiled from further primary literature (a
363 complete list of items is provided in appendix 1). The items were literally translated into German
364 and, where necessary, adapted contextually for the current study. All items were rated on a five-point
365 Likert scale ranging from 1 = Strongly disagree to 5 = Strongly agree. The Presence scale consists

366 of four items, such as ‘I was completely captivated by the virtual environment,’ and originates
367 from the study by Makransky et al. (2017).

368 The Motivation scale (Makransky and Petersen, 2019) comprises five items (e.g., ‘I enjoy working
369 with the topic of titration’), with the word ‘titration’ mentioned in each item to establish context.
370 Similarly, the extrinsic Cognitive Load scale (Andersen and Makransky, 2021) consists of two
371 subscales pertaining to the environment (e.g., ‘The virtual environment was full of irrelevant
372 content’) and control within the virtual space (e.g., ‘The interaction technique used in the simulation
373 was difficult to master’). These subscales measure the demand on cognitive resources consistent with
374 Cognitive Load Theory by Sweller et al. (2011). All three scales (Presence, Motivation, Cognitive
375 Load) demonstrate good reliability with Cronbach’s Alpha values ranging between .8 and .9, as
376 reported by Petersen et al. (2022).

377 2.3.4 Technology Affinity

378 To control for interindividual differences in technology affinity (TA), the short version of the TAEG
379 questionnaire was used (Karrer-Gauß et al., 2024). This scale consists of four items addressing
380 self-assessed domains: competence in using technology (e.g., ‘I am knowledgeable about electronic
381 devices’), enthusiasm for new technologies (e.g., ‘I get excited when a new electronic device comes
382 onto the market’), and the assessment of positive (e.g., ‘Electronic devices make my everyday
383 life easier’) or negative consequences (e.g., ‘Electronic devices lead to mental impoverishment’)
384 associated with using technological applications. Items were rated on a five-point Likert scale. This
385 four items showed a high correlation of .92 in the original study what makes them a suitable control
386 variable for the present study (Karrer-Gauß et al., 2024).

387 2.3.5 System Usability Scale (SUS)

388 To assess subjective usability, we utilized the German translation of the System Usability Scale
389 (SUS) (Brooke, 1996) by Gao et al. (2020). This scale comprises 10 items rated on a five-point Likert
390 scale. For the present study, the items were adapted such that the word ‘product’ was replaced with
391 ‘VR environment’ in all questions. An example item read: ‘I found the various functions in this VR
392 environment were well integrated.’

393 The final SUS score was calculated according to the standard procedure described by Brooke
394 (1996). For items with odd numbers, 1 was subtracted from the score; for items with even numbers,
395 the score was subtracted from 5. The sum of these adjusted scores for an individual was then
396 multiplied by 2.5 to obtain a final score ranging from 0 to 100. The German version of the scale
397 demonstrated good reliability, with Cronbach’s Alpha ranging from $\alpha = .74$ to $\alpha = .88$ (Gao et al.,
398 2020).

399 2.4 Procedure

400 The study was conducted on the Muttentz campus over four days between October 9th and 17th,
401 2024, as part of a mandatory course for the students. Sessions commenced at 7:30 AM, 8:30 AM,
402 and 12:40 PM, with each session lasting approximately 3 hours and 10 minutes on average. During
403 each session, at least two members of the research team were present to administer and support the
404 VR condition, alongside a minimum of four university instructors.

405 Researchers and instructors welcomed the students in the lecture hall and explained the purpose
406 and procedure of the study. Subsequently, participants signed an informed consent form and received

407 an anonymized code for completing the questionnaires. The first questionnaire (t_1 , see Figure 1)
408 collected demographic data, assessed technology affinity, and included the knowledge test.

409 After administering the first questionnaire, students formed pairs, which were then divided into
410 two main groups. Pair formation was initiated by the students themselves without direction from
411 the research team, while assignment to the starting condition was determined by the instructors
412 based on participant numbers. One group was guided to the physical laboratory by the instructors,
413 while the other group received an introduction to the VR headsets. This VR introduction involved
414 demonstrating and explaining the software application's user interface, highlighting the tutorial
415 videos available within the virtual lab, and pointing out the function of various buttons. Researchers
416 assisted students with setting up and using the headsets until they successfully reached the virtual lab
417 and grasped the fundamentals of navigation and interaction within the simulation. Subsequently, the
418 pairs conducted the virtual titration experiment independently, following the instructions provided
419 within the application. The progress of each pair was monitored via a separate Windows application,
420 and researchers intervened only if students were unable to proceed.

421 Concurrently, students assigned to the physical laboratory condition were introduced by the
422 instructors to the equipment and expected procedures. They received a printed version of the
423 experimental instructions and performed the titration independently under instructor supervision.
424 Interventions by the instructors occurred only in response to safety-related issues or when students
425 requested assistance.

426 Upon completion of their first titration (either VR or physical lab), all students returned to the
427 lecture hall to complete the second questionnaire (t_2). The group that had been in the physical
428 laboratory completed the knowledge test again. The group that had experienced the VR condition
429 completed the knowledge test along with the CAMIL and SUS questionnaires. Following this, there
430 was a short break at the students' disposal. Subsequently, the groups switched conditions. After
431 completing their second assigned condition, students filled out the third and final questionnaire (t_3),
432 concluding the session.

433 2.5 Data Analysis

434 The assumptions of the CAMIL model were tested using the Partial Least Squares Structural
435 Equation Modeling (PLS-SEM) approach, as this method, unlike covariance-based SEM models,
436 accommodates smaller sample sizes and does not require normally distributed data (Hair et al.,
437 2019). The analysis and evaluation of the PLS models followed the process described by Hair et al.
438 (2021), and calculations were performed using the SEMinR package (2.3.4, Ray et al., 2024).

439 The evaluation process for the measurement model involved four steps, beginning with the
440 assessment of indicator reliability, followed by internal consistency reliability, convergent validity,
441 and finally, discriminant validity (Hair et al., 2019, p. 76). The evaluation of the structural model
442 comprised five steps, starting with checking for collinearity, evaluating path significance, assessing
443 the model's explanatory and predictive power, and finally, conducting model comparisons when
444 different structural model specifications were considered.

445 To assess the main effects of the test condition (VR or Lab), timepoint, and sequence (the order
446 of VR and Lab), a Linear Mixed-Effects Model (LMM) was calculated, specifying participants
447 as a random effect. The LMM was computed using the lme4 package (1.1.36, Bates et al., 2015).
448 The performance package (0.13.0, Lüdtke et al., 2021) was used to check model assumptions and

449 variance explained, and the effectsize package (1.0.0, Ben-Shachar et al., 2020) was employed to
450 determine the effect sizes of the predictors. The recommendations of Meteyard and Davies (2020)
451 were followed regarding the analysis procedure and the reporting of results.

3 RESULTS

452 Participants were assigned to the different sessions by the module coordinator; consequently, this
453 allocation was not controlled by the research team. To minimize the potential influence of this
454 assignment process, key group characteristics were compared and tested for significant differences.
455 These characteristics were selected as they were measured prior to the first learning condition and
456 are independent of the experimental manipulation. Table 1 displays descriptive statistics for sample
457 size (n), technology affinity (TAEG), baseline test score, age, gender, and educational level, broken
458 down by test group (session). As homogeneity of variances could not be assumed, a Welch's ANOVA
459 was chosen for these comparisons.

460 Table 2 summarizes the results of the Welch's ANOVA tests comparing these variables across
461 the test groups; all p-values were greater than .208. Based on the pre-defined alpha level of .05,
462 there were no significant differences between the groups. The effect sizes (ω^2) for baseline knowledge
463 ($\omega^2 = 0.05$) and age ($\omega^2 = 0.06$) were small (Kirk, 1996), while the effect size for technology affinity
464 ($\omega^2 < 0.001$) was negligible. Additionally, Fisher's exact tests revealed no significant association
465 between session group and gender ($p = .262$) or educational level ($p = .083$). Therefore, it can be
466 assumed that the experimental groups did not differ substantially from one another based on these
467 baseline characteristics.

468 3.1 Knowledge test

469 Table 3 presents the means and standard deviations for the knowledge test scores across all
470 timepoints and conditions. Overall, at baseline (t_1), students achieved an average score of 38.97
471 points (SD = 3.77) out of a maximum possible 52 points, corresponding to 74.9%. This suggests a
472 high level of prior knowledge among the participants.

473 Following the first learning phase (t_2), scores averaged across both conditions increased to 41.64
474 points (SD = 4.31), representing an improvement of 2.67 points over baseline. After the second
475 learning phase (t_3), the average score reached 42.10 points (SD = 4.09). While this still represents
476 an improvement compared to the baseline test, it reflects only a minimal further increase of 0.46
477 points from the preceding test score (t_2). Figure 3 visualizes this pattern, plotting the change scores
478 relative to baseline for all participants across all timepoints, differentiated by starting condition (VR
479 or Lab). A clear increase in scores from t_1 to t_2 is evident, which subsequently levels off between
480 t_2 and t_3 . Notably, this trend was not uniform; considerable variance in test score changes was
481 observed among participants.

482 3.1.1 Psychometric Properties of the Knowledge Test

483 The reliability values for the knowledge test exhibit a specific pattern (cf. Table 4). At the first
484 timepoint t_1 both Cronbach's α and ω are low ($\alpha_{t_1} = .485$ and $\omega_{total_{t_1}} = .563$). However, they
485 increase to their highest values at the second administration ($\alpha_{t_2} = .618$ and $\omega_{total_{t_2}} = .672$), before
486 decreasing again at the third administration ($\alpha_{t_3} = .579$ and $\omega_{total_{t_3}} = .647$), though not below the
487 initial values at t_1 . Simultaneously ω_h increases steadily, reaching its maximum value at the third
488 assessment point ($\omega_{h_{t_1}} = .084$, $\omega_{h_{t_2}} = .196$ and $\omega_{h_{t_3}} = .262$)

489 The initial increase in reliability likely reflects the consolidation of knowledge following the learning
490 task. The continuous increase in ω_h suggests the development of distinct subdimensional skills or
491 knowledge domains, as this measure, unlike α and ω_{total} , accounts for multidimensionality (McNeish,
492 2018).

493 Since the knowledge test likely encompasses multiple facets (e.g., declarative knowledge,
494 mathematical calculations, procedural tasks), the increase in α and ω_{total} – despite the assumption
495 of unidimensionality – provides a useful, albeit cautious, assessment of improved overall reliability.

496 The subsequent decrease in α and ω_{total} at t_3 could be attributed to factors such as fatigue,
497 decreased motivation, or a ceiling effect. Further investigation is required to determine the
498 statistical significance and practical relevance of these fluctuations. As highlighted by Kelley
499 and Pornprasertmanit (2016) and McNeish (2018) α can vary across repeated measures, particularly
500 if the assumption of tau-equivalence is not met.

501 3.1.2 Statistical Analysis of Knowledge Test Scores

502 An overall comparison of mean scores across the three timepoints (t_1, t_2, t_3) showed a significant
503 difference ($F_{Welch}(2, 179.44) = 17.00, p < .001, \omega_p^2 = 0.15, CI_{95\%}[0.07, 1.00], n = 273$). Post-hoc
504 pairwise comparisons indicated significant differences between the first timepoint ($M = 38.97$) and
505 both the second timepoint ($M = 41.64, p < .001$) and third timepoint ($M = 42.10, p < .001$).

506 When analyzing each experimental condition separately, we found no significant increase in
507 test scores for either the lab condition ($t_{Welch}(87.8) = -0.06, p = 0.95, \hat{g}_{Hedges} = -0.01,$
508 $CI_{95\%}[-0.42, 0.40]$) or the VR condition ($t_{Welch}(88.94) = -0.90, p = 0.37, \hat{g}_{Hedges} = -0.19,$
509 $CI_{95\%}[-0.60, 0.22]$). Additionally, a direct comparison between the two conditions at t_2 revealed no
510 statistical difference ($t_{Welch}(88.92) = 1.01, p = 0.32, \hat{g}_{Hedges} = 0.21, CI_{95\%}[-0.20, 0.62], n_{obs} = 91$).
511 This suggests that the overall effect is only attributed to the first condition, regardless of which
512 condition it was, while the second experimental condition showed no further influence on learning
513 outcomes.

514 3.2 SUS Score

515 The overall SUS score for our sample was $M = 71.73(SD = 13.81, n = 91)$, which indicates that
516 the usability of the lab simulation can be interpreted as good or acceptable (Lewis, 2018). The
517 internal consistency was adequate (Cronbach's $\alpha = .817$), consistent with reliability values reported
518 for the German version of the SUS (Gao et al., 2020).

519 3.3 Path analysis of the CAMIL Model

520 The PLS-SEM model (cf. Figure 4) was initially estimated using all available items according
521 to the path-weighting scheme. To determine indicator reliability, the squared factor loadings of
522 all items were compared against the cut-off value of 0.708, although this threshold is not always
523 met in the social sciences, making lower values potentially acceptable (Hair et al., 2021). Following
524 the recommendations of Hair et al. (2022), however, items with squared loadings below 0.4 were
525 excluded from further analysis due to their weak explanatory power for indicator variance. This
526 was the case for the items listed in Table 5, and a subsequent model was estimated excluding these
527 items.

528 The exclusion of a relatively large number of items (8 out of 17) fundamentally questions the
529 transferability of the original scales to the context of VR-based chemistry education. This challenge
530 cannot be solely explained by negatively worded items or translation issues. Particularly concerning
531 the Motivation scale, where three out of five items had to be removed, the resulting construct likely
532 captures only a limited aspect of the intended motivation. This reduction in construct breadth
533 might partially explain why our structural model exhibited low explanatory power and why only
534 the negative relationship between Presence and Cognitive Load emerged as significant. For future
535 research, this suggests a need to adapt and validate the measurement concepts of the CAMIL model
536 and their German translation, specifically for the application domain of VR learning in science
537 laboratory environments.

538 The internal consistency reliability and convergent validity values (cf. Table 6) were acceptable for
539 the restricted model. Furthermore, discriminant validity, assessed using the heterotrait-monotrait
540 (HTMT) method, was established, with all values falling below the threshold of 0.95 (Henseler
541 et al., 2015). Therefore, the measurement model of the restricted model was deemed appropriate for
542 subsequent analysis.

543 To assess the structural model's goodness-of-fit, potential collinearity issues were examined first.
544 Variance Inflation Factor (VIF) values were calculated for the predictors of the endogenous constructs.
545 In our model, this primarily concerns the predictors Cognitive Load (CL) and Motivation (MO)
546 for the Learning Outcome (WIS) construct. The VIF value was 1.019 for both CL and MO,
547 indicating excellent results with no multicollinearity concerns (Becker et al., 2015). Regarding the
548 evaluation of the path coefficients (cf. Table 7), only the negative relationship between Presence
549 (PRE) and Cognitive Load (CL) emerged as statistically significant ($\beta = -0.471$, $t = -6.491$, $p <$
550 $.01$, $CI_{95\%}[-0.630, -0.349]$). Table 8 presents the R^2 values for the endogenous constructs, serving
551 as a measure of the model's explanatory power. The results show that all R^2 values fell below the
552 .250 threshold, which is generally considered weak explanatory power. CL had the highest R^2 at
553 .222, while MO ($R^2 = .045$), PRE ($R^2 = .033$), and WIS ($R^2 = .021$) accounted for very little
554 variance, suggesting they did not make a substantial contribution to the explained variance in this
555 sample.

556 The final two steps for evaluating a PLS-SEM model, as recommended by Hair et al. (2021),
557 were not performed. This decision was based on two reasons: firstly, standard procedures for
558 assessing out-of-sample predictive power (such as PLSpredict) are generally not applicable to models
559 containing higher-order constructs (higher-order composites), CL construct utilized in this study.
560 Secondly, there was insufficient theoretical justification to support the specification of an alternative
561 structural model.

562 In conclusion, the collected data and the resulting structural model demonstrate limited explanatory
563 power. Of the relationships derived from the CAMIL model, only a significant negative influence of
564 Presence on Cognitive Load was confirmed ($\beta = -0.471$, $p < .01$). The direction of this relationship
565 aligns with the findings of Petersen et al. (2022), suggesting an interplay between these constructs
566 and supporting the notion that cognitive load can impair performance in immersive VR environments.
567 Notably, this relationship explained approximately 22% of the variance in Cognitive Load ($R^2 = .222$),
568 indicating that the sense of presence might be an important factor in reducing cognitive burden
569 within the iVR environment.

570 The lack of significant paths from Motivation to Learning Outcome and from Cognitive Load
571 to Learning Outcome contradicts the theoretical propositions of the CAMIL model. This finding
572 might be partially attributable to the aforementioned measurement challenges, but it could also
573 suggest that the learning mechanisms involved in mastering titration within VR differ from those
574 observed for other learning content previously studied in VR contexts. The very low R^2 values for
575 Learning Outcome (.021) and Motivation (.045) suggest that crucial influencing factors, such as
576 prior knowledge, might be absent from the model. Alternatively, the knowledge assessment, with its
577 primary focus on declarative knowledge, may not have adequately captured potential learning gains
578 related to practical laboratory skills.

579 3.4 Linear Mixed Effects Models

580 To further investigate the effects of experimental condition, intervention sequence, and participants,
581 a linear mixed-effects model (LMM) was computed. The outcome variable was the change in test
582 score relative to the initial baseline score (t_1). This approach allowed modeling the change directly
583 without including the baseline score as a covariate. The rationale was that the observed changes
584 were relatively small compared to the initial scores, and including the initial test score as a covariate
585 might have masked the potential effects of the other predictors.

586 Participant was specified as a random effect, estimating a random intercept for each participant.
587 This accounts for the nested data structure (repeated measures within participants) inherent in the
588 study design. Following the recommendation by Barr et al. (2013) to implement a maximal random
589 effects structure, this would ideally involve including random slopes for the effect of Condition (VR
590 vs. Lab) per participant. However, attempting to estimate random slopes for Condition (a factor
591 with only two levels) led to model convergence issues. Therefore, the maximal feasible random
592 effects structure retained only the random intercepts for participants (cf. Eager and Roy, 2017).

593 To determine the best model fit, a baseline model including only the random intercept for
594 participants was estimated first (cf. Table 9). Fixed effects were then added sequentially. The
595 main predictor of interest, experimental Condition, was added as the first fixed effect (Model 1).
596 Subsequently, Timepoint was added in Model 2, and Sequence (VR-first vs. Lab-first) was added in
597 Model 3. All models met the assumptions of linearity, normality of residuals, and homoscedasticity
598 (Meteyard and Davies, 2020).

599 Table 10 presents the model comparison results, indicating no statistically significant improvement
600 in model fits including fixed effects compared to the baseline random-intercept-only model.

601 None of the fixed effects significantly improved model fit or explained additional variance in
602 the outcome variable. It is worth noting, however, that Model 1 (adding Condition) approached
603 statistical significance ($p = 0.06$) and showed a marginal improvement in fit over the baseline
604 model according to the likelihood ratio test ($\chi^2(1) = 3.540$) and AIC comparison ($\Delta AIC = -1.54$),
605 although not BIC ($\Delta BIC = 1.66$).

606 The lack of significant improvement from adding fixed effects, as determined by the model
607 comparisons, is clearly reflected in the variance explained (R^2) by the final model. The variance
608 accounted for by the random effect of participants was substantial (Conditional $R^2 \approx .761$, indicating
609 ~76% of variance attributed to participant differences), whereas the fixed effects collectively explained
610 very little variance (Marginal $R^2 \approx .008$, or 0.8%). Examining the fixed effect for condition revealed a
611 non-significant trend suggesting potentially lower performance in the VR condition compared to the

612 Lab condition ($\beta = -0.526$, $SE = 0.289$, 95%-KI $[-1.09, 0.05]$, $t(df) = -1.817$, $p = .071$, $\eta^2 = 0.036$).
613 However, due to the lack of statistical significance, no firm conclusion can be drawn regarding this
614 difference (see Table 11). Taken together, these results indicate that inter-individual differences in
615 learning performance accounted for the vast majority of the observed variance, while neither the
616 specific intervention condition (VR vs. Lab) nor the sequence of interventions had a statistically
617 detectable impact on the measured learning outcomes.

4 DISCUSSION

618 The principal goal of this study was to shed some light on the learning efficiency of an iVR training
619 sequence compared to the same activity in a conventional laboratory setting in an actual mandatory
620 first-semester Life-Sciences module. In line with recent research (Johnson-Glenberg et al., 2023; Liu
621 et al., 2023), we found no improvement in the learning outcome between the two conditions, as shown
622 in model fit comparisons when *Condition* (VR vs. Lab) was entered in our mixed-effects analysis,
623 knowledge gains were statistically equivalent across modalities. Simple pair-wise comparisons at
624 the second timepoint corroborated this null effect. Moreover, learning outcomes were unaffected
625 by the order in which students encountered the modalities (VR→Lab vs. Lab→VR), indicating
626 that neither the environment itself nor its sequencing confers an advantage. These results suggest
627 that iVR can serve as a pedagogically equivalent, be potentially more scalable and complementary
628 to traditional laboratories. An insight that should inform future investment decisions in teaching
629 infrastructure and could guide instructional designers when creating new learning activities.

630 4.1 Implications for Learning in iVR

631 Our first hypothesis (H1) posited that students' technology affinity would enhance their sense
632 of presence in the simulated laboratory. The path analysis did reveal a positive standardized
633 coefficient, but this effect was not statistically significant. Consequently, within the present sample
634 and measurement precision we cannot claim that technology affinity reliably boosts presence. Because
635 the confidence interval still encompasses small-to-moderate effects, follow-up studies with larger
636 samples or refined instruments remain warranted before ruling out a meaningful relationship.

637 In line with the CAMIL model and recent meta-analytic evidence (Concannon et al., 2019; Liu
638 et al., 2023; Parong and Mayer, 2018), a stronger sense of presence was associated with lower
639 extraneous cognitive load. Put differently, the more students felt "present" the virtual lab, the
640 fewer cognitive resources were consumed by interface demands or distractions from the simulated
641 environment, leaving greater capacity for task-relevant processing. H2 was therefor confirmed. This
642 finding highlights the importance of designing iVR environments that actively foster presence while
643 minimizing non-essential elements that could compete for learners' cognitive resources. Bøg Petersen
644 et al. (2022) found that learners in VR environments often retain less knowledge compared to
645 more passive learning methods, due to the cognitive load imposed by the simulation. Similarly,
646 Parong and Mayer (2021a) refer to cognitive or affective distraction in such contexts. Another
647 cognitive mechanism is reported by Oje et al. (2025), who point out, that the use of leading cues and
648 guidance in VR, what we did in our virtual laboratory, can lower the extraneous CL and facilitate
649 the concentration on the learning task. Our result may indicate that the sense of presence could
650 serve as an important predictor of reduced cognitive load—an effect that should, in turn, support
651 improved learning outcomes.

652 Surprisingly, two further CAMIL paths did not materialize: presence was not a reliable predictor
653 of motivation, H3, and neither cognitive load nor motivation predicted knowledge-test performance
654 (H4). Measurement attenuation likely played a role here; several motivation items fell below reliability
655 thresholds and were removed, narrowing construct coverage and weakening statistical power. While
656 the negative association between motivation and test score can likely be attributed to issues with
657 the measurement instrument, the lack of a relationship between cognitive load and test score is more
658 difficult to explain. As previously discussed, lower cognitive load should generally be associated with
659 better learning outcomes (Parong and Mayer, 2021a; Sweller et al., 2011). However, this connection
660 could not be confirmed in our model. An explanation might be a ceiling effect in the test scores,
661 which could obscure potential effects due to limited variance in the outcome variable. Another
662 consideration is the design of the test instrument itself, which may have focused too heavily on
663 declarative knowledge and insufficiently captured the procedural components of the learning process
664 (Hamilton et al., 2021).

665 The structural model's fixed effects explained only $\approx 0.8\%$ of the variance in test scores (marginal
666 $R^2 \approx .008$), whereas random intercepts for participants accounted for roughly 76% (conditional
667 $R^2 \approx .76$). Such dominance of individual differences suggests that prior knowledge or test-taking skills,
668 rather than the experimental manipulations, drove most outcome variability. This interpretation
669 dovetails with the high baseline score (about 75% correct) and the very modest absolute gain
670 between the second and third knowledge tests (approximately 0.5 points). Together, the pattern is
671 consistent with ceiling effects, limited test sensitivity to procedural gains, and diminishing returns
672 once the titration protocol has been practiced in either modality.

673 4.2 Considerations about sequencing and technologies

674 The fifth hypothesis (H5) examined whether differences exist between the instructional methods
675 employed and their sequencing. Interpretation of this linear mixed-effects model must again consider
676 that the fixed effects — condition, timepoint, and sequence — accounted for only 0.8% of the
677 variance explained by the model. As observed in the path model, the majority of the variance
678 ($\sim 76\%$) was attributable to the random effect of the individual participant. This finding reinforces
679 the conclusion that a substantial proportion of the observed effects can be explained by inherent
680 individual characteristics. Moreover, by employing a second statistical approach, we were able to
681 confirm that there was no significant difference between the two learning activities. The sequence
682 of the activities likewise showed no significant influence, further supporting our assertion that
683 similar learning gains can be achieved when comparable processes or actions are implemented across
684 activities. These findings corroborate previous meta-analyses (Concannon et al., 2019; Liu et al.,
685 2023), which similarly reported no, or only minimal, differences between iVR and other media or
686 instructional methods. The main difference with these studies is that we conducted this study in
687 an actual course for a grounding skill in this curriculum. We can therefore conclude, that in this
688 very realistic setting, the simulation was on the same level than the laboratory learning activity
689 regarding this specific test and his results.

690 4.3 Limitations

691 This study has a pilot character, as it represents the first implementation of this learning activity
692 within the module. The scales used to measure motivation, presence, and cognitive load were not
693 validated in advance and only demonstrate satisfactory validity. This likely influenced the statistical
694 results. During the VR condition, some groups received more support and guidance when they

695 encountered difficulties during the practice phase, which likely introduced unequal treatment between
696 groups. In addition, group assignments were not fully randomized. The influence of collaboration
697 during the simulation phase was neither measured nor explicitly guided. Furthermore, participants
698 in the VR condition often had to wait longer for the laboratory condition participants to return
699 before completing the knowledge test. This delay may have resulted in memory decay effects, as the
700 lab group had the learning content fresher in mind. Conversely, motivational issues may have arisen
701 from the idle waiting time. Finally, the third questionnaire round always took place at the end of
702 the school block, either shortly before lunch or before going home, which—alongside the waiting
703 period for the VR group—may have contributed to fatigue or a tendency to complete this final
704 phase quickly and with reduced attention.

705 **4.4 Practical Implications and future research**

706 The demonstrated equivalence in learning outcomes between the simulation and the real laboratory
707 underscores the potential of immersive technology as a viable means to enrich student education with
708 an additional dimension. However, as [Radianti et al. \(2020\)](#) emphasizes, its implementation must be
709 holistic and grounded in learning theory. Assuming, as [Conrad et al. \(2024\)](#) suggest, that there is
710 no significant difference between various instructional technologies and media, and considering that
711 iVR can be particularly effective for procedural knowledge acquisition ([Hamilton et al., 2021](#)), that
712 there is compelling evidence for transfer from VR to real-world workplaces ([Bøg Petersen et al.,
713 2022](#)), and that—according to our findings—neither sequence nor condition had a decisive impact,
714 the focus must now shift to the appropriate pedagogical integration of this technology. If VR is to
715 be treated as an equivalent pedagogical tool, curricula must be intentionally designed so that this
716 additional modality contributes meaningfully to the development of targeted competencies, attitudes
717 or skills. Instructional design models such as ADDIE ([Branch, 2009](#)) or Constructive Alignment
718 ([Biggs and Tang, 2011](#)) may provide valuable frameworks in this regard. These approaches can be
719 combined with evidence-based design and pedagogical principles in VR as reported by [Oje et al.
720 \(2025\)](#), with the aim to create a model for designing learning activities in VR.

721 Future research should also engage more deeply with the assessment of procedural learning in
722 VR (see also [Hamilton et al., 2021](#)) and explore how the resulting data can support personalized
723 and adaptive learning pathways. Which data points could serve as indicators of successful learning?
724 How can gaps in students' knowledge be derived from this data and addressed directly? These
725 questions are likely to become increasingly relevant, particularly as AI models are integrated either
726 in the background to guide and support learning processes or more actively as agents or learning
727 companions within the simulation. It will also be essential to examine how students accept these
728 mechanisms and how they affect cognitive and affective processes.

729 From an economic perspective, it is important to evaluate whether the acquisition, development,
730 and maintenance of the technological infrastructure required to deliver such simulations offer
731 advantages over traditional real-world learning environments. Equally critical is the initial and
732 ongoing training of educators in the effective use and pedagogical application of these technologies.

733 **4.5 Conclusion**

734 Contemporary educational design increasingly emphasizes competency development and acquisition
735 in the construction of curricula and training programs. Our work aligns with this evolution, and
736 in light of current AI-driven disruptions, the importance of concrete actions in education and

737 professional development will undoubtedly increase. In connection with this trend, personal, concrete
738 experiential learning opportunities are extremely valuable. However, resources are limited regarding
739 physical spaces, facilities, and personnel when implementing real-world environments for educational
740 purposes.

741 The application of simulations (iVR especially) is of significant interest in this context. In
742 our specific case, we demonstrated that the implementation of iVR compared to real laboratory
743 experiences in a curricular learning activity and subsequent assessment led to identical outcomes.
744 Based on other studies, we can reasonably assume this finding may be generalizable, suggesting that
745 the application of simulations in educational offerings should undergo more thorough examination,
746 and thoughtful implementation of this technology will certainly become increasingly important in
747 the future.

748 This approach offers substantial advantages for students, enabling them to engage with learning
749 experiences not only at educational institutions and in real environments but also asynchronously
750 and across various devices. Therefore, the utilization of simulations presents a viable opportunity
751 when seeking to establish educational equity or addressing varying learning rates, offering significant
752 benefits by potentially allowing students to repeat such learning experiences multiple times as
753 needed.

754 It is worth noting, however, that further research is required to establish a robust and reliable
755 framework for designing these activities in VR (and in near future also in mixed reality settings).
756 The development of virtual learning experiences necessitates evidence-based guidance, and their
757 integration into educational organizations requires a thoughtful approach that involves instructors in
758 the process. Only through such collaborative efforts can we fully realize the potential of immersive
759 technologies in educational contexts.

CONFLICT OF INTEREST

760 The author declares no conflict of interest. The study was conducted as part of the author's Master
761 thesis in the applied psychology program at the FHNW.

ACKNOWLEDGEMENTS

762 I would like to thank Dr. Oliver Christ for providing me with the opportunity to participate in
763 this study and to write this thesis. His guidance, feedback, and support throughout the writing
764 process were invaluable. I would also like to thank the teaching team from the Life Sciences Bachelor
765 Course for their collaboration in creating the questionnaire, their assistance in testing it, and their
766 constructive feedback. Their contribution was essential during the conduct of the study, as they
767 adapted and organized their modules to accommodate the iVR condition and integrated it seamlessly
768 into their regular coursework.

769 Use of Generative AI Technologies: During the writing of this work, various generative AI models
770 were used to assist with translation of German text to English, grammar and syntax correction,
771 brainstorming ideas and arguments, and providing feedback on text sections. No AI models were
772 used to generate original content or write text in the first instance. All writing represents my original
773 intellectual work, with AI tools used solely for correction, translation, and feedback purposes.
774 Portions of the R code used for statistical analysis and results reporting were generated by AI

775 following my specific instructions. All calculations and reported results are based exclusively on
776 the actual research data processed through R code, not on AI model outputs. The following AI
777 models were utilized: OpenAI GPT-4o, GPT-4.1, GPT-3o; Anthropic Claude 3.5 Sonnet, Claude
778 3.7 Sonnet, Claude Sonnet 4, Claude Opus 4; Google Gemini 1.5 Flash, Gemini 1.5 Pro, Gemini 2.0
779 Flash, Gemini 2.0 Pro; Google NotebookLM, GitHub Copilot.

DATA AVAILABILITY

780 The data that support the findings of this study are available from the corresponding author upon
781 reasonable request.

BIBLIOGRAPHY

- 782 Agustian, H. Y., Finne, L. T., Jørgensen, J. T., Pedersen, M. I., Christiansen, F. V., Gammelgaard,
783 B., et al. (2022). Learning outcomes of university chemistry teaching in laboratories: A systematic
784 review of empirical literature. *Review of Education* 10, e3360, doi: [10.1002/rev3.3360](https://doi.org/10.1002/rev3.3360)
- 785 Andersen, M. S. and Makransky, G. (2021). The validation and further development of a
786 multidimensional cognitive load scale for virtual environments. *Journal of Computer Assisted*
787 *Learning* 37, 183–196, doi: [10.1111/jcal.12478](https://doi.org/10.1111/jcal.12478)
- 788 Bailenson, J. (2018). *Experience on demand: what virtual reality is, how it works, and what it can*
789 *do* (New York, NY London: W.W. Norton & Company)
- 790 Barr, D. J., Levy, R., Scheepers, C., and Tily, H. J. (2013). Random effects structure for
791 confirmatory hypothesis testing: Keep it maximal. *Journal of Memory and Language* 68, 255–278,
792 doi: [10.1016/j.jml.2012.11.001](https://doi.org/10.1016/j.jml.2012.11.001)
- 793 Barrett, A., Pack, A., Guo, Y., and Wang, N. (2023). Technology acceptance model and multi-
794 user virtual reality learning environments for chinese language education. *Interactive Learning*
795 *Environments* 31, 1665–1682, doi: [10.1080/10494820.2020.1855209](https://doi.org/10.1080/10494820.2020.1855209)
- 796 Bates, D., Mächler, M., Bolker, B., and Walker, S. (2015). Fitting linear mixed-effects models using
797 lme4. *Journal of Statistical Software* 67, 1–48, doi: [10.18637/jss.v067.i01](https://doi.org/10.18637/jss.v067.i01)
- 798 Becker, J.-M., Ringle, C. M., Sarstedt, M., and Völckner, F. (2015). How collinearity affects mixture
799 regression results. *Marketing Letters* 26, 643–659, doi: [10.1007/s11002-014-9299-9](https://doi.org/10.1007/s11002-014-9299-9)
- 800 Ben-Shachar, M. S., Lüdtke, D., and Makowski, D. (2020). effectsize: Estimation of effect size indices
801 and standardized parameters. *Journal of Open Source Software* 5, 2815, doi: [10.21105/joss.02815](https://doi.org/10.21105/joss.02815)
- 802 Bennett, S. W. and O’Neale, K. (1998). Skills development and practical work in chemistry.
803 *University Chemistry Education* 2
- 804 Biggs, J. B. and Tang, C. S.-k. (2011). *Teaching for quality learning at university: what the student*
805 *does*. SRHE and Open University Press Imprint (Maidenhead: McGraw-Hill/Society for Research
806 into Higher Education/Open University Press), 4th edition edn.
- 807 Bohne, T., Heine, I., Gurerk, O., Rieger, C., Kemmer, L., and Y. Cao, L. (2021). Perception
808 engineering learning with virtual reality. *IEEE Transactions on Learning Technologies* 14, 500–514,
809 doi: [10.1109/TLT.2021.3107407](https://doi.org/10.1109/TLT.2021.3107407)
- 810 Branch, R. M. (2009). *Instructional Design: The ADDIE Approach* (Boston, MA: Springer US), doi:
811 [10.1007/978-0-387-09506-6](https://doi.org/10.1007/978-0-387-09506-6)
- 812 Bretz, S. L. (2019). Evidence for the importance of laboratory courses. *Journal of Chemical*
813 *Education* 96, 193–195, doi: [10.1021/acs.jchemed.8b00874](https://doi.org/10.1021/acs.jchemed.8b00874)

- 814 Brooke, J. (1996). Sus: A 'quick and dirty' usability scale. In *Usability Evaluation in Industry*, eds.
815 P. W. Jordan, B. Thomas, I. L. McClelland, and B. Weerdmeester (London: Taylor & Francis).
816 189–194, doi: [10.1201/9781498710411](https://doi.org/10.1201/9781498710411)
- 817 Bøg Petersen, G., Klingenberg, S., and Makransky, G. (2022). Pipetting in virtual reality can predict
818 real-life pipetting performance. *Technology, Mind, and Behavior* 3, doi: [10.1037/tmb0000076](https://doi.org/10.1037/tmb0000076)
- 819 Campbell, C. D., Midson, M. O., Bergstrom Mann, P. E., Cahill, S. T., Green, N. J. B., Harris,
820 M. T., et al. (2022). Developing a skills-based practical chemistry programme: an integrated, spiral
821 curriculum approach. *Chemistry Teacher International* 4, 243–257, doi: [10.1515/cti-2022-0003](https://doi.org/10.1515/cti-2022-0003)
- 822 Chirico, A. and Gaggioli, A. (2019). When virtual feels real: Comparing emotional responses and
823 presence in virtual and natural environments. *Cyberpsychology, Behavior, and Social Networking*
824 22, 220–226, doi: [10.1089/cyber.2018.0393](https://doi.org/10.1089/cyber.2018.0393)
- 825 Concannon, B. J., Esmail, S., and Roduta Roberts, M. (2019). Head-mounted display virtual
826 reality in post-secondary education and skill training. *Frontiers in Education* 4, 80, doi:
827 [10.3389/feduc.2019.00080](https://doi.org/10.3389/feduc.2019.00080)
- 828 Conrad, M., Kablitz, D., and Schumann, S. (2024). Learning effectiveness of immersive virtual
829 reality in education and training: A systematic review of findings. *Computers & Education: X*
830 *Reality* 4, 100053, doi: [10.1016/j.cexr.2024.100053](https://doi.org/10.1016/j.cexr.2024.100053)
- 831 Cooper, N., Millela, F., Cant, I., White, M. D., and Meyer, G. (2021). Transfer of training—virtual
832 reality training with augmented multisensory cues improves user experience during training and
833 task performance in the real world. *PLOS ONE* 16, e0248225, doi: [10.1371/journal.pone.0248225](https://doi.org/10.1371/journal.pone.0248225)
- 834 di Lanzo, J. A., Valentine, A., Sohel, F., Yapp, A. Y. T., Muparadzi, K. C., and Abdelmalek, M.
835 (2020). A review of the uses of virtual reality in engineering education. *Computer Applications in*
836 *Engineering Education* 28, 748–763, doi: [10.1002/cae.22243](https://doi.org/10.1002/cae.22243)
- 837 Dunnagan, C. L., Dannenberg, D. A., Cuales, M. P., Earnest, A. D., Gurnsey, R. M., and Gallardo-
838 Williams, M. T. (2020). Production and evaluation of a realistic immersive virtual reality organic
839 chemistry laboratory experience: Infrared spectroscopy. *Journal of Chemical Education* 97,
840 258–262, doi: [10.1021/acs.jchemed.9b00705](https://doi.org/10.1021/acs.jchemed.9b00705)
- 841 Eager, C. and Roy, J. (2017). Mixed effects models are sometimes terrible, doi:
842 [10.48550/arXiv.1701.04858](https://doi.org/10.48550/arXiv.1701.04858)
- 843 European Commission and PwC (2020). *Skills for industry curriculum guidelines 4.0: future proof*
844 *education and training for manufacturing in Europe : final report*. Tech. rep., Publications Office,
845 LU
- 846 Finne, L. T., Gammelgaard, B., and Christiansen, F. V. (2022). When the lab work disappears:
847 Students' perception of laboratory teaching for quality learning. *Journal of Chemical Education*
848 99, 1766–1774, doi: [10.1021/acs.jchemed.1c01113](https://doi.org/10.1021/acs.jchemed.1c01113)
- 849 Gao, M., Kortum, P., and Oswald, F. L. (2020). Multi-language toolkit for the system
850 usability scale. *International Journal of Human-Computer Interaction* 36, 1883–1901, doi:
851 [10.1080/10447318.2020.1801173](https://doi.org/10.1080/10447318.2020.1801173)
- 852 Guo, H., Ma, F., and Zhou, Z. (2025). Validation of technology acceptance model for virtual reality
853 usage in collaborative learning to enhance learner performance. *Innovations in Education and*
854 *Teaching International* 62, 429–443, doi: [10.1080/14703297.2024.2307994](https://doi.org/10.1080/14703297.2024.2307994)
- 855 Hair, J. F., Black, W. C., Babin, B. J., and Anderson, R. E. (2019). *Multivariate data analysis*
856 (Andover, Hampshire: Cengage), eighth edition edn.

- 857 Hair, J. F., Hult, G. T. M., Ringle, C. M., and Sarstedt, M. (2022). *A primer on partial least squares*
858 *structural equation modeling (PLS-SEM)* (Thousand Oaks: SAGE Publications, Incorporated),
859 third edition edn.
- 860 Hair, J. F., Hult, G. T. M., Ringle, C. M., Sarstedt, M., Danks, N. P., and Ray, S. (2021).
861 *Partial Least Squares Structural Equation Modeling (PLS-SEM) Using R: A Workbook*. Classroom
862 Companion: Business (Cham: Springer International Publishing), doi: [10.1007/978-3-030-80519-7](https://doi.org/10.1007/978-3-030-80519-7)
- 863 Hamilton, D., McKechnie, J., Edgerton, E., and Wilson, C. (2021). Immersive virtual reality as a
864 pedagogical tool in education: a systematic literature review of quantitative learning outcomes and
865 experimental design. *Journal of Computers in Education* 8, 1–32, doi: [10.1007/s40692-020-00169-2](https://doi.org/10.1007/s40692-020-00169-2)
- 866 Henseler, J., Ringle, C. M., and Sarstedt, M. (2015). A new criterion for assessing discriminant
867 validity in variance-based structural equation modeling. *Journal of the Academy of Marketing*
868 *Science* 43, 115–135, doi: [10.1007/s11747-014-0403-8](https://doi.org/10.1007/s11747-014-0403-8)
- 869 Hyde, J. (2025). A new perspective on chemistry foundation level students laboratory skill
870 development using reciprocal peer-teaching, laboratory simulations, and practical skills portfolio
871 (psp) during covid-19 and post-pandemic in 2024. *Journal of Chemical Education* 102, 984–1003,
872 doi: [10.1021/acs.jchemed.4c01124](https://doi.org/10.1021/acs.jchemed.4c01124)
- 873 Johnson-Glenberg, M. C., Yu, C. S. P., Liu, F., Amador, C., Bao, Y., Yu, S., et al. (2023). Embodied
874 mixed reality with passive haptics in stem education: randomized control study with chemistry
875 titration. *Frontiers in Virtual Reality* 4, 1047833, doi: [10.3389/frvir.2023.1047833](https://doi.org/10.3389/frvir.2023.1047833)
- 876 Johnstone, A. H. (1993). The development of chemistry teaching: A changing response to changing
877 demand. *Journal of Chemical Education* 70, 701, doi: [10.1021/ed070p701](https://doi.org/10.1021/ed070p701)
- 878 Kaplan, A. D., Cruit, J., Endsley, M., Beers, S. M., Sawyer, B. D., and Hancock, P. A. (2021). The
879 effects of virtual reality, augmented reality, and mixed reality as training enhancement methods:
880 A meta-analysis. *Human Factors: The Journal of the Human Factors and Ergonomics Society* 63,
881 706–726, doi: [10.1177/0018720820904229](https://doi.org/10.1177/0018720820904229)
- 882 Karrer-Gauß, K., Glaser, C., Clemens, C., and Bruder, C. (2009). Technikaffinität erfassen – der
883 fragebogen ta-eg. *Der Mensch im Mittelpunkt technischer Systeme* 8, 196–201
- 884 Karrer-Gauß, K., Roesler, E., and Siebert, F. W. (2024). Neuauflage des taeg fragebogens:
885 Technikaffinität valide und multidimensional mit einer kurz- oder langversion erfassen. *Zeitschrift*
886 *für Arbeitswissenschaft* 78, 387–406, doi: [10.1007/s41449-024-00427-4](https://doi.org/10.1007/s41449-024-00427-4)
- 887 Kelley, K. and Pornprasertmanit, S. (2016). Confidence intervals for population reliability coefficients:
888 Evaluation of methods, recommendations, and software for composite measures. *Psychological*
889 *Methods* 21, 69–92, doi: [10.1037/a0040086](https://doi.org/10.1037/a0040086)
- 890 Kirk, R. E. (1996). Practical significance: A concept whose time has come. *Educational and*
891 *Psychological Measurement* 56, 746–759, doi: [10.1177/0013164496056005002](https://doi.org/10.1177/0013164496056005002)
- 892 Levac, D. E., Huber, M. E., and Sternad, D. (2019). Learning and transfer of complex motor skills
893 in virtual reality: a perspective review. *Journal of NeuroEngineering and Rehabilitation* 16, 121,
894 doi: [10.1186/s12984-019-0587-8](https://doi.org/10.1186/s12984-019-0587-8)
- 895 Lewis, J. R. (2018). The system usability scale: Past, present, and future. *International Journal of*
896 *Human–Computer Interaction* 34, 577–590, doi: [10.1080/10447318.2018.1455307](https://doi.org/10.1080/10447318.2018.1455307)
- 897 Liu, J. Y. W., Yin, Y.-H., Kor, P. P. K., Cheung, D. S. K., Zhao, I. Y., Wang, S., et al. (2023).
898 The effects of immersive virtual reality applications on enhancing the learning outcomes of
899 undergraduate health care students: Systematic review with meta-synthesis. *Journal of Medical*
900 *Internet Research* 25, e39989, doi: [10.2196/39989](https://doi.org/10.2196/39989)

- 901 Lüdecke, D., Ben-Shachar, M., Patil, I., Waggoner, P., and Makowski, D. (2021). performance: An r
902 package for assessment, comparison and testing of statistical models. *Journal of Open Source*
903 *Software* 6, 3139, doi: [10.21105/joss.03139](https://doi.org/10.21105/joss.03139)
- 904 Makransky, G., Andreasen, N. K., Baceviciute, S., and Mayer, R. E. (2021). Immersive virtual
905 reality increases liking but not learning with a science simulation and generative learning strategies
906 promote learning in immersive virtual reality. *Journal of Educational Psychology* 113, 719–735,
907 doi: [10.1037/edu0000473](https://doi.org/10.1037/edu0000473)
- 908 Makransky, G. and Lilleholt, L. (2018). A structural equation modeling investigation of the emotional
909 value of immersive virtual reality in education. *Educational Technology Research and Development*
910 66, 1141–1164, doi: [10.1007/s11423-018-9581-2](https://doi.org/10.1007/s11423-018-9581-2)
- 911 Makransky, G., Lilleholt, L., and Aaby, A. (2017). Development and validation of the multimodal
912 presence scale for virtual reality environments: A confirmatory factor analysis and item response
913 theory approach. *Computers in Human Behavior* 72, 276–285, doi: [10.1016/j.chb.2017.02.066](https://doi.org/10.1016/j.chb.2017.02.066)
- 914 Makransky, G. and Petersen, G. (2019). Investigating the process of learning with desktop virtual
915 reality: A structural equation modeling approach. *Computers & Education* 134, 15–30, doi:
916 [10.1016/j.compedu.2019.02.002](https://doi.org/10.1016/j.compedu.2019.02.002)
- 917 Makransky, G. and Petersen, G. B. (2021). The cognitive affective model of immersive learning
918 (camil): a theoretical research-based model of learning in immersive virtual reality. *Educational*
919 *Psychology Review* 33, 937–958, doi: [10.1007/s10648-020-09586-2](https://doi.org/10.1007/s10648-020-09586-2)
- 920 McNeish, D. (2018). Thanks coefficient alpha, we'll take it from here. *Psychological Methods* 23,
921 412–433, doi: [10.1037/met0000144](https://doi.org/10.1037/met0000144)
- 922 Meteyard, L. and Davies, R. A. (2020). Best practice guidance for linear mixed-effects
923 models in psychological science. *Journal of Memory and Language* 112, 104092, doi:
924 [10.1016/j.jml.2020.104092](https://doi.org/10.1016/j.jml.2020.104092)
- 925 Meyer-Drawe, K. (2003). Lernen als erfahrung. *Zeitschrift für Erziehungswissenschaft* 6, 505–514,
926 doi: [10.1007/s11618-003-0054-x](https://doi.org/10.1007/s11618-003-0054-x)
- 927 Milgram, P., Takemura, H., Utsumi, A., and Kishino, F. (1995). Augmented reality: a class of
928 displays on the reality-virtuality continuum. In *Telemanipulator and Telepresence Technologies*,
929 ed. H. Das (Boston, MA: SPIE), vol. 2351, 282–292, doi: [10.1117/12.197321](https://doi.org/10.1117/12.197321)
- 930 Oje, A. V., Hunsu, N. J., and Fiorella, L. (2025). A systematic review of evidence-based design and
931 pedagogical principles in educational virtual reality environments. *Educational Research Review*
932 47, 100676, doi: [10.1016/j.edurev.2025.100676](https://doi.org/10.1016/j.edurev.2025.100676)
- 933 Parong, J. and Mayer, R. (2018). Learning science in immersive virtual reality. *Journal of*
934 *Educational Psychology* doi: [10.1037/EDU0000241](https://doi.org/10.1037/EDU0000241)
- 935 Parong, J. and Mayer, R. E. (2021a). Cognitive and affective processes for learning science
936 in immersive virtual reality. *Journal of Computer Assisted Learning* 37, 226–241, doi:
937 [10.1111/jcal.12482](https://doi.org/10.1111/jcal.12482)
- 938 Parong, J. and Mayer, R. E. (2021b). Learning about history in immersive virtual reality: does
939 immersion facilitate learning? *Educational Technology Research and Development* 69, 1433–1451,
940 doi: [10.1007/s11423-021-09999-y](https://doi.org/10.1007/s11423-021-09999-y)
- 941 Peeters, H., Habig, S., and Fechner, S. (2023). Does augmented reality help to understand chemical
942 phenomena during hands-on experiments?—implications for cognitive load and learning. *Multimodal*
943 *Technologies and Interaction* 7, 9, doi: [10.3390/mti7020009](https://doi.org/10.3390/mti7020009)
- 944 Petersen, G., Petkakis, G., and Makransky, G. (2022). A study of how immersion and interactivity
945 drive vr learning. *Computers & Education* 179, 104429, doi: [10.1016/j.compedu.2021.104429](https://doi.org/10.1016/j.compedu.2021.104429)

- 946 Radianti, J., Majchrzak, T. A., Fromm, J., and Wohlgenannt, I. (2020). A systematic review of
947 immersive virtual reality applications for higher education: Design elements, lessons learned, and
948 research agenda. *Computers & Education* 147, 103778, doi: [10.1016/j.compedu.2019.103778](https://doi.org/10.1016/j.compedu.2019.103778)
- 949 Rahman, A. A. (2020). Tracing the evolution of transfer of training: A review article. *Annals of*
950 *Social Sciences & Management studies* 5, doi: [10.19080/ASM.2020.05.555668](https://doi.org/10.19080/ASM.2020.05.555668)
- 951 Ray, S., Danks, N. P., and Valdez, A. (2024). seminar: Building and estimating structural equation
952 models
- 953 Sheppard, K. (2006). High school students' understanding of titrations and related acid-base
954 phenomena. *Chem. Educ. Res. Pract.* 7, 32–45, doi: [10.1039/B5RP90014J](https://doi.org/10.1039/B5RP90014J)
- 955 Sweller, J., Ayres, P., and Kalyuga, S. (2011). *Cognitive Load Theory* (New York, NY: Springer
956 New York), doi: [10.1007/978-1-4419-8126-4](https://doi.org/10.1007/978-1-4419-8126-4)
- 957 Tee, N. Y. K., Gan, H. S., Li, J., Cheong, B. H.-P., Tan, H. Y., Liew, O. W., et al. (2018).
958 Developing and demonstrating an augmented reality colorimetric titration tool. *Journal of*
959 *Chemical Education* 95, 393–399, doi: [10.1021/acs.jchemed.7b00618](https://doi.org/10.1021/acs.jchemed.7b00618)
- 960 Thisgaard, M. and Makransky, G. (2017). Virtual learning simulations in high school: Effects on
961 cognitive and non-cognitive outcomes and implications on the development of stem academic and
962 career choice. *Frontiers in Psychology* 8, 805, doi: [10.3389/fpsyg.2017.00805](https://doi.org/10.3389/fpsyg.2017.00805)
- 963 Tusher, H. M., Mallam, S., and Nazir, S. (2024). A systematic review of virtual reality features for
964 skill training. *Technology, Knowledge and Learning* 29, 843–878, doi: [10.1007/s10758-023-09713-2](https://doi.org/10.1007/s10758-023-09713-2)
- 965 UNESCO (2024). *Higher education: figures at a glance*. Tech. rep., UNESCO
- 966 Wolter, S. C., Albiez, J., Cattaneo, M. A., Denzler, S., Diem, A., Lüthi, S., et al. (2023).
967 *Bildungsbericht Schweiz 2023* (Aarau: Schweizerische Koordinationsstelle für Bildungsforschung
968 (SKBF))
- 969 Won, M., Ungu, D. A. K., Matovu, H., Treagust, D. F., Tsai, C.-C., Park, J., et al. (2023).
970 Diverse approaches to learning with immersive virtual reality identified from a systematic review.
971 *Computers & Education* 195, 104701, doi: [10.1016/j.compedu.2022.104701](https://doi.org/10.1016/j.compedu.2022.104701)
- 972 Zierer, K. and Seel, N. M. (2012). General didactics and instructional design: eyes like twins a
973 transatlantic dialogue about similarities and differences, about the past and the future of two
974 sciences of learning and teaching. *SpringerPlus* 1, 15, doi: [10.1186/2193-1801-1-15](https://doi.org/10.1186/2193-1801-1-15)

FIGURE CAPTIONS

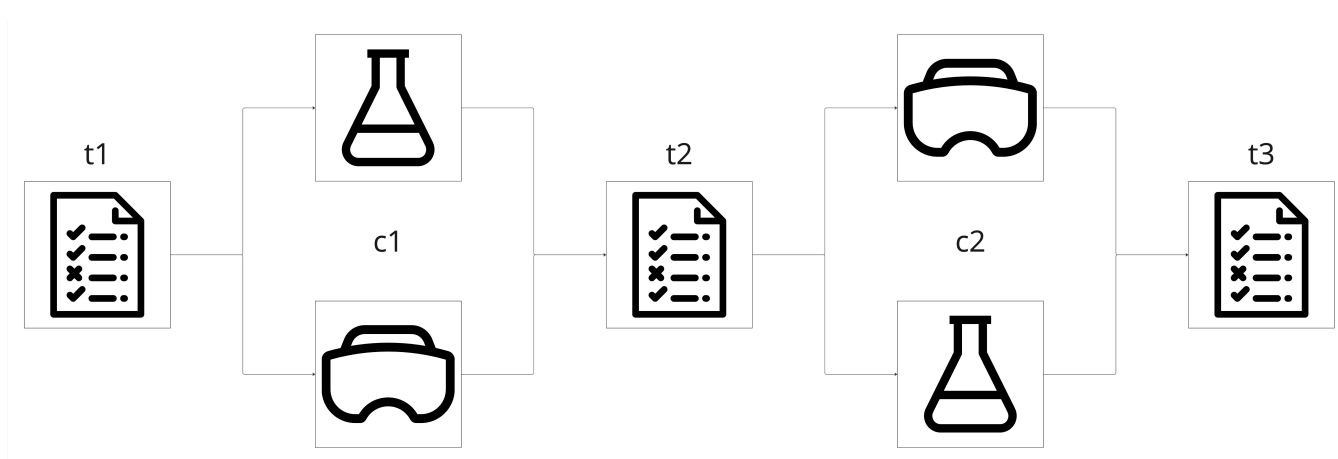


Figure 1. Schematic diagram of the study procedure

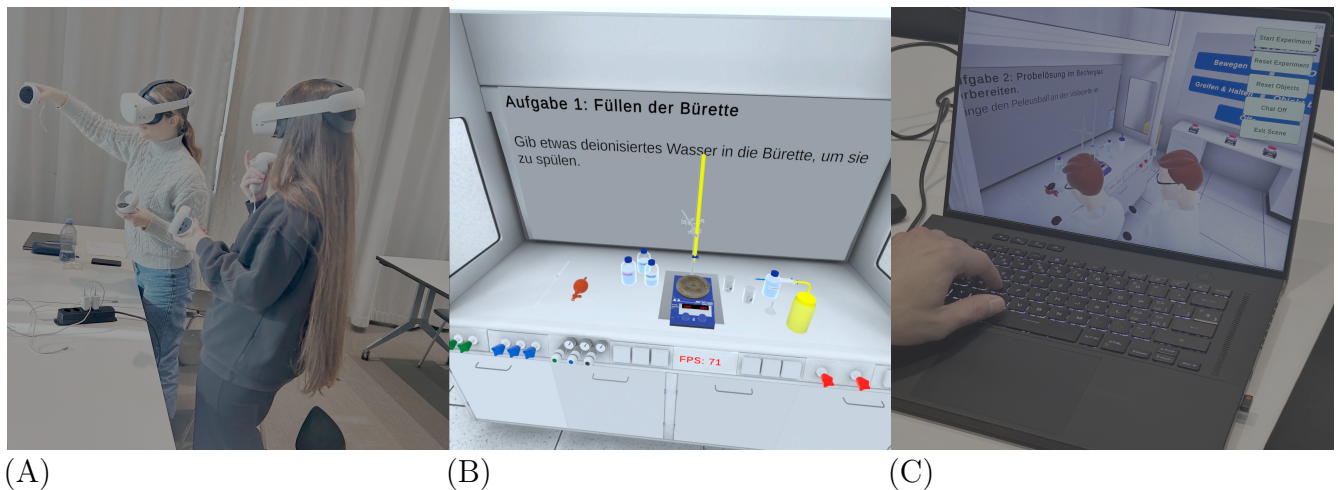


Figure 2. Figures of the VR condition; (A) Students working in iVR, (B) iVR working place, (C) iVR Supervision

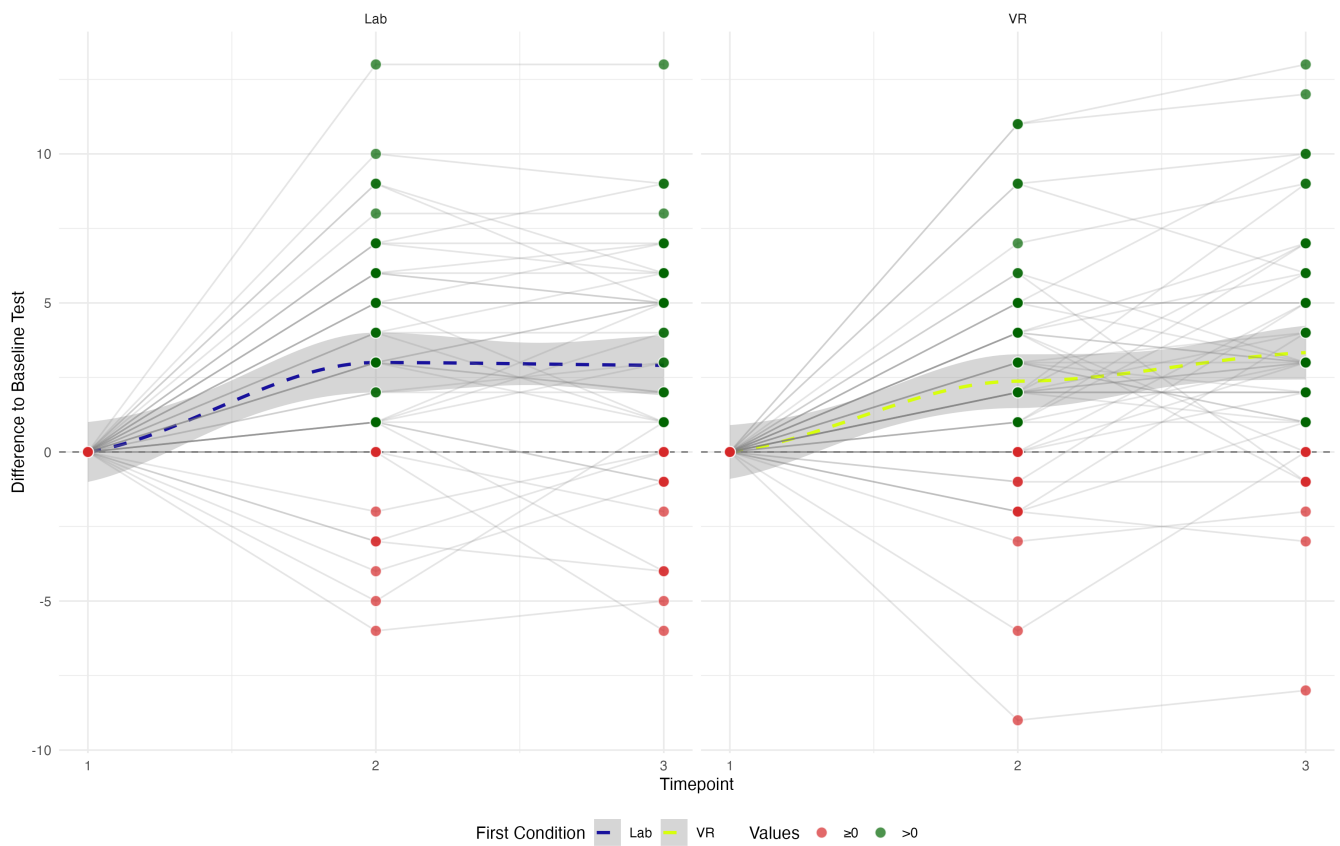


Figure 3. Knowledge test score trends by starting condition and timepoint

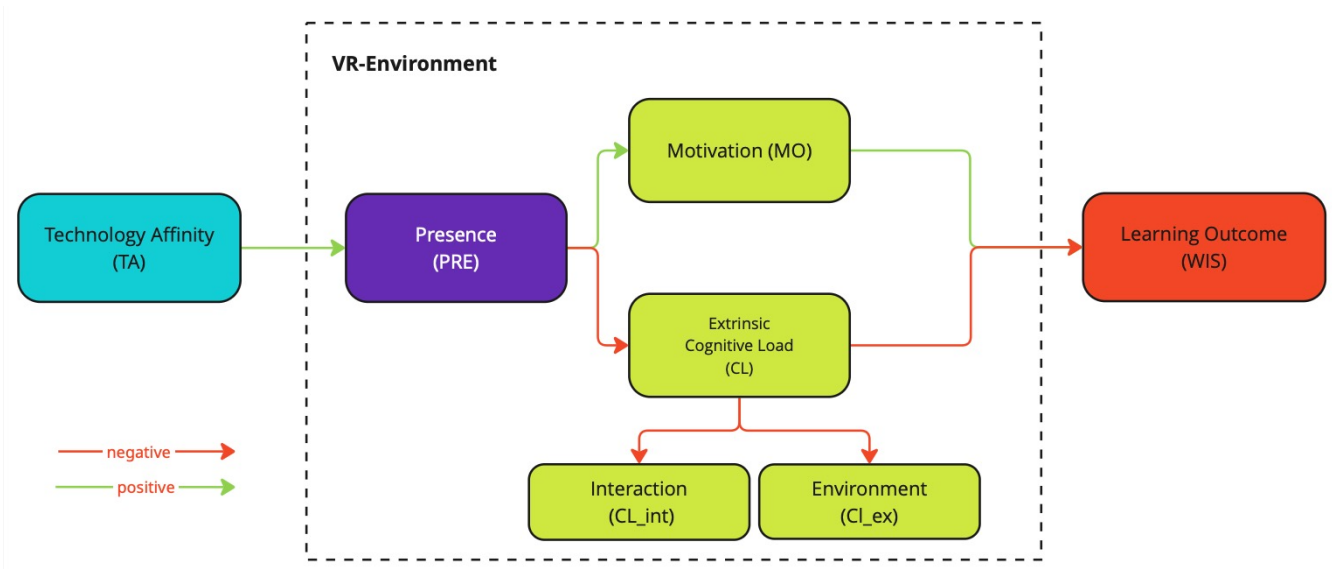


Figure 4. PLS-SEM of the CAMIL model (Makransky and Petersen, 2021)

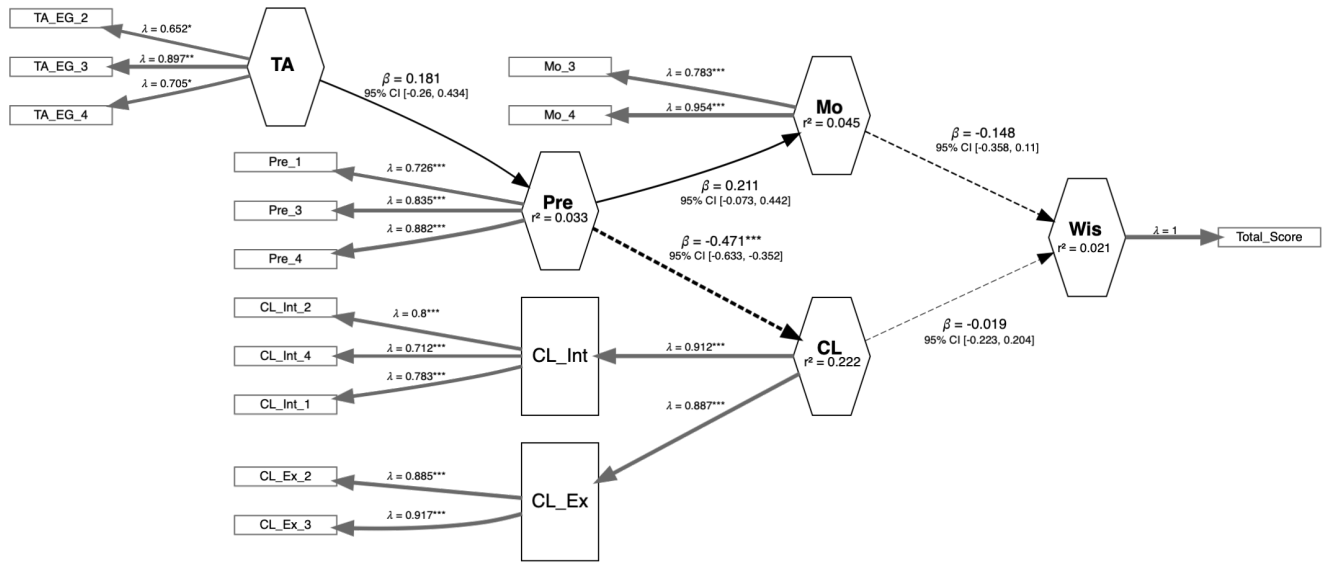


Figure 5. PLS-SEM with path coefficients

TABLES

Table 1 Descriptive statistics of the experimental groups

Session	n	TAEG		Baseline		Age		Sex		
		Mean	Sd	Mean	Sd	Mean	Sd	f	m	d
1	12	3.27	0.76	40.00	4.24	22.17	3.21	9	3	
2	20	3.52	0.44	39.00	3.78	20.80	1.74	11	9	
3	17	3.59	0.63	40.24	2.66	20.88	2.18	11	5	1
4	10	3.45	0.57	37.00	4.00	22.60	2.41	5	5	
5	18	3.18	0.70	38.17	3.65	22.06	3.49	9	9	
6	14	3.55	0.56	38.93	4.25	23.29	6.51	12	2	

Table 2 One-way Welch’s ANOVA results comparing session groups

Variable	F	df1	df2	p	ω^2	n
Baseline	1.41	5	35.38	.245	0.05	91
Age	1.52	5	34.53	.208	0.06	91
TA_EG	0.95	5	35.36	.459	0.00	91

Table 3 Results of the knowledge tests

Condition	Timepoint	n	Mean	SD
Baseline	1	91	38.97	3.77
Lab	2	43	42.12	4.12
Lab	3	48	42.17	4.10
VR	2	48	41.21	4.48
VR	3	43	42.02	4.11

Table 4 Reliability values for test items

Scale	Numberof items	CronbachAlpha	OmegaTotal	OmegaHierarchical
TA_EG	4	.650	.713	.427
Mo	5	.822	.860	.643
Pre	4	.700	.765	.353
CL	8	.789	.827	.508
t1	13	.422	.507	.014
t2	13	.632	.678	.052
t3	13	.577	.628	.307

Table 5 Items with factor loadings with values < 0.4

	Squared Loadings
TA_EG_1	0.374
Mo_1	0.336
Mo_2	0.069
Mo_5	0.076
Pre_2	0.364
CL_Ex_1	0.241
CL_Ex_4	0.262
CL_Int_3	0.220

Table 6 Internal consistency reliability and convergent validity metrics

	α	ρ_C	ρ_A	AVE
TA	0.648	0.800	0.790	0.576
Pre	0.758	0.857	0.817	0.667
CL	0.765	0.895	0.773	0.809
Mo	0.719	0.864	1.001	0.762
CL_Ex	0.769	0.896	0.782	0.812
CL_Int	0.653	0.809	0.667	0.587

Table 7 Path Coefficients of the Bootstrapped Structural Model ($n=10,000$)

	Original Estimate	Bootstrap Mean	SD	T Stat.	p	Confidence Interval	
						2.5%	97.5%
TA → Pre	0.181	0.185	0.180	1.008	.316	-0.255	0.431
Pre → CL	-0.471	-0.497	0.072	-6.538	<.001	-0.630	-0.351
Pre → Mo	0.211	0.221	0.130	1.632	.106	-0.073	0.449
CL → Wis	-0.019	-0.016	0.109	-0.172	.864	-0.222	0.201
Mo → Wis	-0.148	-0.150	0.115	-1.285	.202	-0.354	0.105

Table 8 R²-Values of the Bootstrapped Structural Model ($n=10,000$)

	Pre	CL	Mo	Wis
R ²	0.033	0.222	0.045	0.021
adj. R ²	0.022	0.213	0.034	-0.001
TA	0.181			
Pre		-0.471	0.211	
CL				-0.019
Mo				-0.148

Table 9 Model specifications with fit indices

Model	Fixed Effects	AIC	BIC	Log-Likelihood	df
Null	None	947.83	957.44	-470.91	3
Condition	Condition	946.29	959.10	-469.14	4
Timepoint	Timepoint + Condition	946.03	962.05	-468.01	5
Sequence	Timepoint + Condition + Sequence	948.01	967.24	-468.01	6

N total obs = 182; N Subjects = 91

Table 10 Likelihood Ratio Tests Between Nested Models

Model	Compared to	df	χ^2	p
Condition	Null	1	3.54	.060
Timepoint	Condition	1	2.26	.133
Sequence	Timepoint	1	0.02	.898

Table 11 Comprehensive Overview of Fixed and Random Effects in the Model

	β	SE	Fixed Effects 95% CI	t	p	η^2
Intercept	2.135	0.937	0.24 - 3.93	2.279	.024	
Timepoint (1)	0.433	0.289	-0.14 - 1.00	1.496	.136	0.025
Condition (VR)	-0.526	0.289	-1.09 - 0.05	-1.817	.071	0.036
Sequence (VR first)	-0.099	0.781	-1.63 - 1.44	-0.127	.899	0.000

Note: The p-values for the Fixed Effects were calculated using Satterthwaite approximations. The confidence intervals were calculated using the bootstrap method based on 10000 simulations. Model equation: Base_Diff_Score \sim Timepoint + Condition + Sequence + (1 | User_Code)