

# Counting your chickens before they hatch: improvements in an untreated chronic pain population, beyond regression to the mean and the placebo effect

Monica Sean<sup>1,2,5</sup>, Alexia Coulombe-Lévêque<sup>1,4,6</sup>, William Nadeau<sup>1</sup>, Anne-Catherine Charest<sup>1</sup>, Marylie Martel<sup>1,5</sup>, Guillaume Léonard<sup>1,4,6</sup>, Pascal Tétreault<sup>\*1,2,3,5</sup>

\*Corresponding Author: [Pascal.Tetreault@USherbrooke.ca](mailto:Pascal.Tetreault@USherbrooke.ca)

<sup>1</sup>Université de Sherbrooke, Faculty of medicine and health sciences, Sherbrooke, Québec, Canada

<sup>2</sup>Department of Anesthesiology, Sherbrooke, Québec, Canada

<sup>3</sup>Department of Nuclear Medicine and Radiobiology

<sup>4</sup>School of Rehabilitation, Sherbrooke, Québec, Canada

<sup>5</sup>Centre de recherche du CHUS, Sherbrooke, Québec, Canada

<sup>6</sup>Research Centre on Aging, Sherbrooke, Québec Canada.

## Abstract

**Background and aims:** Isolating the effect of an intervention from the natural course and fluctuations of a condition is a challenge in any clinical trial, particularly in the field of pain. Regression to the mean (RTM), wherein extreme scores are more likely to be followed by more average scores, may explain some of those observed fluctuations. However, while this phenomenon is relatively well-known, its effect on outcome measures is rarely quantified, and often only evoked as a potential confound. In this paper, we describe and quantify such symptom fluctuations in a chronic pain population in the absence of treatment, and compare the relative stability of various self-reported outcome measures in untreated chronic low back pain (CLBP) patients and healthy controls (HC).

**Methods:** Twenty-three untreated CLBP patients and 25 HC took part in this observational study, wherein they were asked to complete an array of commonly used questionnaires in pain studies (including the Pain Catastrophizing Scale [PCS], State and Trait Anxiety Inventory [STAI], Central Sensitization Inventory [CSI], Pain Disability Index [PDI], Brief Pain Inventory [BPI] etc.) during each of three visits (V1, V2, V3) at 2-month intervals. Scores at V1 were classified into three subgroups (extremely high, normal and extremely low), based on z-scores. The average delta ( $\Delta=V2-V1$ ) was calculated for each subgroup, for each questionnaire, to describe the evolution of scores over time. This analysis was repeated with the data for V2 and V3.

**Results:** High initial scores were likely to be followed by more average scores; for instance, the “extremely high” subgroup for the PCS (a reputedly ‘stable’ questionnaire) had an average decrease of 12/52 from V1 to V2. Participants with “average” initial scores tended to show a small decrease over time, and participants with “extremely low” initial scores tended to remain stable. However, while the pattern of fluctuation in the three subgroups was similar across questionnaires, the magnitude of these fluctuations varied greatly. The STAI and CSI were the most stable questionnaires of all, with even

47 the “extremely high” subgroup showing little or no improvement over time. The least  
48 stable questionnaires were the PCS, PDI and BPI.

49

50 **Discussion and conclusion:** These pain trajectories in untreated patients cannot be  
51 attributable to RTM alone because of their asymmetry, nor to the placebo effect as they  
52 occurred in the absence of any intervention. We propose that the observed improvements  
53 could be the result of an Effect of Care, wherein participants had meaningful  
54 improvements simply from taking part in a study and talking about their pain to  
55 benevolent research staff, despite the absence of active or sham treatment. These findings  
56 have important clinical ramifications. Beyond simply raising a flag as to the existence  
57 (and significance) of Effect of Care, we provide a questionnaire- specific baseline of  
58 expected fluctuations based on initial score, against which researchers can compare  
59 results from clinical trials when trying to isolate the effect of their intervention.

60

61

## 62 **Keywords**

63 Regression To the Mean (RTM)

64 Effect of Care (EC)

65 Chronic Pain (CP)

66 Pain Questionnaires

67 Outcomes measures

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

87

88

89

90

91

92

93 **Introduction**

94  
95 Chronic pain conditions are highly prevalent and debilitating; as such, many clinical trials  
96 are trying to identify treatment approaches<sup>1-7</sup>. However, because pain is a highly  
97 subjective and variable phenomenon, it is famously difficult to measure accurately – even  
98 more so when it comes to measuring *changes* in pain levels. Indeed, pain levels fluctuate  
99 naturally, as they are affected by a wide array of biopsychosocial factors such as sleep,  
100 mood, expectations, beliefs, etc.<sup>8-11</sup>. These factors, like pain, are also often subjective  
101 and variable, and difficult to assess – and their impact on pain is similarly challenging to  
102 measure with precision<sup>12</sup>.

103  
104 Clinical trials do their best to quantify the effects of their interventions using the most  
105 valid and reliable questionnaires at their disposal<sup>13,14</sup>; unfortunately, it remains difficult  
106 to isolate the effect of an intervention from the natural course and fluctuations of the  
107 condition<sup>15</sup>. This is further complicated by a relatively well-known phenomenon:  
108 regression to the mean (RTM)<sup>16</sup>.

109  
110 RTM is *not* another biopsychosocial factor that *influences* pain levels: it is a *statistical*  
111 *concept* that can – in part – describe, explain, and predict those fluctuations. RTM is  
112 based on probability distributions, and states that extreme scores are likely to be followed  
113 by less extremes scores that are closer to the individual’s own sampling mean<sup>16</sup>. As such,  
114 RTM for pain symptoms depends on two factors: 1) the variability in symptom severity  
115 for a given subject (i.e., whether this subject’s symptoms tend to be very constant over  
116 time, or whether they fluctuate wildly from one day to the next); and 2) the ‘extreme-  
117 ness’ of the subject’s symptoms at baseline, compared to their own sampling distribution  
118 (i.e., how severe their symptoms happened to be on the day of measurement, compared to  
119 their average symptom severity)<sup>15</sup>.

120  
121 RTM is a widely known concept, and it is often mentioned in the discussion section of  
122 clinical trial reports, as a possible alternative explanation for observed changes in  
123 outcomes over time. However, RTM has rarely been the primary focus of a clinical  
124 publication with a chronic pain population. Nevertheless, quantifying RTM – and

125 identifying outcome measures that are intrinsically more susceptible to show RTM – is  
126 theoretically possible. For example, on any given questionnaire, RTM would predict that  
127 an extreme score (high or low) is likely to be followed by a more average score,  
128 regardless of treatment effect; in other words, a high score is expected to be followed by  
129 a lower score (and vice versa), while a more average score is expected to remain  
130 relatively stable. Complex conditions such as chronic pain are affected by a number of  
131 biopsychosocial factors which also have intrinsic variability. Different questionnaires  
132 measure these different factors in different ways, such that it is possible that two  
133 questionnaires show different RTM for the same subject over the same time period. In  
134 other words, a questionnaire measuring a comparatively more fluctuating factor will be  
135 more susceptible to RTM.

136

137 Quantifying RTM in the questionnaires most often used to assess treatment efficacy in  
138 the chronic pain population is of clinical significance for two reasons: 1) It may guide the  
139 choice of outcome measures selected at the time of study design; and 2) It may improve  
140 result interpretation, helping to differentiate treatment effect from RTM.

141

142 The study design required for such RTM assessment requires that patients with chronic  
143 pain be assessed using a large array of validated, commonly used questionnaires, at  
144 different time points (at least twice, ideally more), with no concomitant intervention  
145 taking place outside of usual care. Our team had such an observational study taking place  
146 to assess changes in brain structure and functional activity over time in patients with  
147 chronic low-back pain (CLBP) and healthy controls (HC). We were therefore able to  
148 conduct the RTM analysis presented below as part of this study.

149

150 The objectives of this analysis were: 1) to describe and quantify the natural trajectory of  
151 questionnaire scores over time, based on initial scores, with a subgoal of determining  
152 whether the observed fluctuations were compatible with RTM; and 2) to evaluate and  
153 compare the stability of each questionnaire over time, in 23 untreated CLBP and 25  
154 healthy controls.

155

156

## 157 **Methods**

158 All participants provided written informed consent for their participation into the study.  
159 Ethics approval was granted from the institutional review board of the Centre intégré  
160 universitaire de santé et de services sociaux de l'Estrie - Centre hospitalier universitaire  
161 de Sherbrooke (CIUSSS de l'Estrie - CHUS), Sherbrooke, Canada (file number: 2021-  
162 3861). The trial has been registered on Open Science Framework (OSF), under the name  
163 "Pilot project on brain and lower back imaging of chronic pain" (registration DOI:  
164 <https://doi.org/10.17605/OSF.IO/P2Z6Y>).

165

### 166 Participants

167 Participants were recruited using posters at the CIUSSS de l'Estrie-CHUS, Facebook ads,  
168 and by word of mouth. Twenty-seven CLBP patients and 25 HC aged 18 to 75 years old  
169 took part in this study. HC were matched with CLBP patients for sex and age.

170 Specific inclusion criteria for the CLBP were: 1) low back pain ( $\geq 6$  months) with or  
171 without pain radiating to the legs or radiating to the neck; 2) average pain intensity of  $\geq$   
172 3/10 in the 24-hour period before the initial visit; 3) pain primarily localized in the lower  
173 back; 4) no history of invasive or aggressive treatment to manage their pain (e.g.  
174 corticosteroid infiltration, strong doses of opioids or antidepressants). Specific exclusion  
175 criteria for HC were: 1) history of chronic pain; 2) pain at the time of testing; 3) an  
176 outstanding painful episode within 3 months of enrollment in the study.

177 Exclusion criteria for the two populations included: 1) neurological, cardiovascular, or  
178 pulmonary disorders; 2) comorbid pain syndrome (i.e., fibromyalgia, osteoarthritis,  
179 irritable bowel syndrome, migraine etc.); 3) history of surgical intervention in the back;  
180 4) used of opioids, antidepressants, anticonvulsants, or psychostimulants; 5) a  
181 corticosteroid infiltration within the past year; 6) pregnancy (current or planned during

182 the course of the study); 7) inability to read or understand French ; 8) contra-indication to  
183 Magnetic Resonance Imaging (MRI).

#### 184 Study design

185 The study had an observational longitudinal design. All participants attended three  
186 sessions at the Centre de recherche du CHUS, where they completed several  
187 questionnaires (reported, analyzed and discussed in the present paper), and provided a  
188 saliva sample and underwent brain and lumbar MRI (as part of the larger study). These  
189 sessions (V1, V2 and V3) took place at two-month intervals.

190

#### 191 Questionnaires

192 CLBP patients completed eight questionnaires: 1) the Pain Catastrophizing Scale (PCS),  
193 2) Pain Disability Index (PDI), 3) Brief Pain Inventory (BPI); 4) Pain DETECT; 5) Pain  
194 Outcomes Questionnaire (POQ), 6) State-Trait Anxiety Inventory (STAI/S-T), 7) McGill  
195 Pain Questionnaire (MPQ), and 8) Central Sensitization Inventory (short form) (CSI).

196

197 HC completed only the PCS and the STAI/S-T (two questionnaires applicable to a  
198 healthy population).

199

200 All questionnaires were completed online using the platform “Research Electronic Data  
201 Capture” (REDCap), and are presented in **Table 1**.

202

203 To avoid fatigue caused by filling out multiple questionnaires, the PDI (for the CLBP  
204 participants) and the PCS (for all participants) were completed at home, one week before  
205 each visit. These two questionnaires were chosen because they make no measure of  
206 "now" (unlike the other questionnaires), instead measuring more general variables that  
207 are unlikely to change over the course of a few days. The remaining questionnaires were  
208 completed at the Centre de recherche du CHUS. Questionnaires were always completed  
209 in the same order, as listed below.

210

211 *Pain Catastrophizing Scale*

212 The Pain Catastrophizing Scale (PCS) is based on three components of pain  
213 catastrophizing: helplessness, magnification, and rumination<sup>13</sup>, and consists of 13 items  
214 rated on a five-point likert-scale (“not at all” to “all the time”) with a score ranging from  
215 0 to 52 (where 52 corresponds to high levels of catastrophizing). The French validated  
216 version was used<sup>17</sup>.

217

#### 218 *The Pain Disability Index*

219 The Pain Disability Index (PDI) is a questionnaire based on different areas of day living  
220 activities such as, home, social, recreational, occupational, sexual, self-care and life  
221 support activities, consisting of seven items rated on a numerical scale (0 = no disability  
222 to 10= worst disability)<sup>18</sup>. PDI total score ranges from 0-70. The French validated version  
223 was used<sup>19</sup>.

224

#### 225 *The Brief Pain Inventory*

226 The Brief Pain Inventory (BPI) short form consists of two subscales, assessing 1) pain  
227 severity (BPIs) and 2) pain interference with function (BPIi)<sup>14</sup>. The BPIs assesses current  
228 pain intensity, and average, worst and least pain in the last 24 hours. These 4 items are  
229 rated on a numerical rating scale (0 = no pain; 10= worst pain imaginable), for a total  
230 score between 0 and 40. The BPIi assesses different functional components such as  
231 mood, sleep, ability to walk, etc., through 7 items rated on a numerical scale (0 = no  
232 interference; 10=complete interference), for a total score between 0 and 70. The French  
233 validated version was used<sup>20</sup>.

234

#### 235 *The PainDETECT*

236 The PainDETECT (PD) evaluates the presence of neuropathic pain components in  
237 patients with back-pain, such as burning sensation, electric shocks etc.<sup>21</sup>. The PD  
238 comprises 7 items related to neuropathic symptoms rated on a 6-point Likert Scale (0=not  
239 at all, 5=very strongly), to which a pain behavior pattern score (-1 to 1) and a radiation  
240 score (0 or 2) are added, for a total score between 0 and 38.

241 The PD also includes 3 pain-intensity scores (current pain intensity, average pain in the  
242 last 4 weeks, and worst pain in the last 4 weeks), which we averaged into a single score

243 called PDs, ranging from 0 to 10 (0 = no pain; 10 = worst pain imaginable). The French  
244 validated version was used<sup>21</sup>.

245

#### 246 *The Pain Outcomes Questionnaire*

247 The Pain Outcomes Questionnaire (POQ) is based on a wide array of components, such  
248 as pain, mobility, activities of daily living, vitality, negative affect and fear<sup>22</sup>. The POQ  
249 comprises 19 items evaluated on a scale from 0 to 10, for a total score ranging from 0  
250 (least symptoms) to 190 (most severe symptoms). There was no validated French  
251 translation available, so we used in-house translation for this questionnaire.

252

#### 253 *The State-Trait Anxiety Inventory*

254 The State-Trait Anxiety Inventory (STAI-S/T) consists in two subscales assessing 1)  
255 State (i.e., current) anxiety, and 2) Trait (i.e., general) anxiety respectively<sup>23</sup>. Each  
256 subscale comprises 20 items rated on a 4-point Likert scale, for a total score between 20  
257 (very low anxiety) to 80 (very high anxiety). The French validated version was used<sup>24</sup>.

258

#### 259 *The McGill Pain Questionnaire*

260 The McGill Pain Questionnaire (MPQ) short form evaluates the sensory-affective  
261 components of pain using 15 items evaluated on a 4-point Likert Scale (0=no pain;  
262 3=severe pain), for a total score ranging from 0 to 45<sup>25</sup>.

263 The MPQ also includes a separate question assessing the average pain intensity over the  
264 previous week, using a 100-point visual analogue scale (left anchor: “no pain”; right  
265 anchor: “worst possible pain”). This score was called MPQi and was analyzed separately  
266 from the rest of the MPQ. The French validated version was used<sup>26</sup>.

267

#### 268 *The Central Sensitization Inventory*

269 The Central Sensitization Inventory (CSI) short form assesses symptoms suggesting the  
270 presence of central sensitization or central sensitivity syndromes, using 25 items rated  
271 on a five-point Likert-scale (0=“never” ; 4= “always”)<sup>27</sup>, with total scores ranging from 0  
272 to 100. The French validated version was used<sup>28</sup>.

273



274

275

## 276 Data Processing and Statistical Analysis

277

### 278 *Group attribution*

279 To first describe the behavior of “extreme” vs “normal” scores, it was necessary to  
280 establish a criterion to differentiate “extreme” and “normal” scores. This was done by  
281 transforming initial raw scores into studentized scores (i.e., z-scores) for each  
282 questionnaire. Multiple z-score thresholds were tested, and a threshold of  $|z| > 0.5$  was  
283 found to yield the most similar number of participants across the 3 subgroups (“extremely  
284 high”, “normal”, “extremely low”). As such, scores with  $|z| > 0.5$  (i.e., scores that were  
285 more than half a standard deviation above or below the group average) were considered  
286 “extreme”, whereas scores with  $|z| < 0.5$  (i.e., scores within half a standard deviation of  
287 the group average) were considered “normal”. An exploratory analysis with various  
288 thresholds revealed that, regardless of the threshold used, a similar pattern emerged from  
289 our results. All results obtained using the different thresholds tested ( $|z| > 0.66$ ;  $|z| > 0,8$   
290 and  $|z| > 1$ ) are included in **supplementary materials**.

291

292 Scores at V1 were thus classified as a) extremely high, b) normal, or c) extremely low.  
293 This was done independently for each questionnaire, such that a given participant could  
294 be in the “extremely high” subgroup for one questionnaire, but in the “normal” subgroup  
295 for another questionnaire. Next, the delta between V1 and V2 was calculated by  
296 subtracting the score at V1 from the score at V2 ( $\Delta=V2-V1$ ), such that a positive delta  
297 corresponds to a score increase (i.e., worsening of the condition), and a negative delta  
298 corresponds to a score decrease (i.e., improvement of the condition). The same analysis  
299 was conducted between V2 and V3: scores at V2 were again classified as a) extremely  
300 high, b) normal, or c) extremely low, and the delta between V2 and V3 was calculated by  
301 subtracting the score at V2 from the score at V3 ( $\Delta=V3-V2$ ). Thus, for each  
302 questionnaire, two calculations were performed (V2-V1 and V2-V3).

303

### 304 *Standardization across questionnaires*

305 To facilitate the comparison between questionnaires, which use various scales, all raw  
306 scores were reported on a scale from 0 to 100. Fluctuations larger than 10 percentage  
307 points between visits were considered clinically meaningful, and fluctuations of 5  
308 percentage points or less were considered random noise. Fluctuations between 5 and 10  
309 percentage points ( $5 \leq |\Delta| \leq 10$ ), while of debatable clinical relevance, were still considered  
310 likely enough to denote an effect to warrant being reported. The use of such standardized  
311 thresholds, as opposed to the Minimal Detectable Change (MDC) specific to each  
312 questionnaire, was favored because it allowed for direct comparisons between  
313 questionnaires; the advantages and drawbacks of this methodological choice are  
314 highlighted in the Discussion.

315

#### 316 *Average delta scores*

317 Once participants were divided into the three subgroups (based on their initial scores),  
318 average delta scores were calculated for each subgroup within each questionnaire. This  
319 yielded a measure of the average evolution over time for each subgroup.

320

#### 321 *Fluctuation scores*

322 In addition to *average* delta scores, “fluctuation scores” were calculated for each  
323 subgroup within each questionnaire by averaging the *absolute value* of delta scores for  
324 the given subgroup. This yielded a measure of the magnitudes of fluctuations at play,  
325 regardless of their direction. This measure was particularly informative in cases where  
326 both large decreases and large increases had taken place. For example, a questionnaire  
327 could have an average score for all participants of 27/100 at V1 and of 28/100 at V2,  
328 yielding an *average* delta score of only 1/100 – seemingly very stable over time.  
329 However, such ‘stability’ could actually be the result of large increases in some  
330 participants and large decreases in others, cancelling each other out. If that were the case,  
331 while the *average* delta would be close to 0, the *fluctuation* score would be large, thus  
332 more accurately representing the variability of scores over time for that particular  
333 questionnaire.

334

#### 335 *Inferential statistics*

336 In line with the exploratory nature of this paper and given the small sample size,  
337 emphasis was placed on descriptive statistics.

338

339

## 340 **Results**

### 341 Participants

342 Fifty-two participants (25 HC and 27 CLBP) were recruited in the study. Three CLBP  
343 participants dropped out after the first visit (unexpected pregnancy [n=1], discomfort  
344 during MRI [n=1], scheduling conflicts [n=1]), and one dropped out after the second visit  
345 (move to a different city [n=1]), such that 23 CLBP completed the entire study and were  
346 included in the analysis. There were no dropouts among the HC.

347

348 The 25 HC (11 women, 12 men) were aged  $44 \pm 15$  years old, and the 23 CLBP  
349 participants (15 women, 11 men) were aged  $40 \pm 14$  years old. Other sociodemographic  
350 characteristics of the sample are presented in **Table 2**. All participants complied with the  
351 instructions relating to medication/treatment throughout the study, namely, that they were  
352 to avoid any treatment other than over-the-counter medication and their usual, non-  
353 invasive rehabilitation treatments. This allowed us to evaluate the natural course of the  
354 condition during the period of the study.

355

356 Participants completed each questionnaire 3 times throughout the study, with 2 months  
357 between each visit. The average scores for each questionnaire at each visit are presented  
358 in **Table 3a and 3b** for the two populations.

359

360 The evolution over time of ‘extremely high’ scorers, ‘extremely low’ scorers, and  
361 ‘normal’ scorers in the CLBP sample (see *Statistical Analysis* for a detailed description of  
362 the analysis method) is shown in Table 4, represented in 3 different ways (Table 4a, 4b,  
363 and 4c). Table 4a shows the *average* evolution for each subgroup. For example, on the  
364 PCS, the CLBP participants in the ‘extremely high’ subgroup at V1 (i.e., participants  
365 with a z-score  $>0.5$ ) had, on average, a reduction of 22 percentage points in their scores at

366 V2 – corresponding to a reduction of 12 points on the PCS scale (range: 0-52). Table 4b  
367 shows the same data but with *individual* delta scores as opposed to *average* delta scores  
368 for each subgroup. For example, for the PCS, the top left cell shows all individual delta  
369 scores (from V1 to V2) for the participants classified in the ‘extremely high’ subgroups at  
370 V1: two participants had a reduction in score ~30 percentage point, and four participants  
371 had a reduction in score ~20 percentage points. These individual scores, averaged  
372 together, yield the average score (-22) reported in table 4a. Table 4c shows the  
373 ‘fluctuation scores’, corresponding to the average of the *absolute value* of the deltas for  
374 each subgroup (in the case of the PCS scores from V1 to V2, for the ‘extremely high’  
375 subgroup, this fluctuation score is equal to the mean delta score because all deltas were  
376 negative; however, for the ‘normal’ subgroup, the fluctuation is score is substantially  
377 larger than the average delta score, because of the presence of both positive and negative  
378 delta scores).

379

380 As stated in the methods, deltas smaller than or equal to 5 percentage points were  
381 considered random noise and most likely not denoting a real change, while deltas larger  
382 10 percentage point were considered clinically meaningful. Deltas between 5 and 10,  
383 while arguably not clinically meaningful, were still considered likely to be more than  
384 simple noise, and therefore to be worth reporting.

385

386

387 Average evolution over time as a function of initial score

388

389 *Average evolution of ‘extremely high’ scores*

390 As shown in the top rows of Table 4a and 4b, participants with an initial ‘extremely high’  
391 score at V1 tended to show a reduction in score at V2, and those with an ‘extremely high’  
392 score at V2 similarly tended to show a reduction in score at V3. Indeed, the *average*  
393 deltas from V1 to V2 and from V2 to V3 were mostly negative in that subgroup (table 4a,  
394 top row) and most *individual* deltas were negative (table 4b, top row). This trend for high  
395 scores to be followed by lower scores is expected, as RTM predicts that extreme scores

396 will be followed by more normal scores, which in the case of extremely *high* scores  
397 means that the subsequent score should be *lower*.

398

#### 399 *Average evolution of ‘normal’ scores*

400 The evolution of ‘normal’ scorers is shown in the second rows of Table 4a, 4b and 4c. As  
401 shown in table 4a, *average* scores tended to remain stable or to slightly decrease over  
402 time. Indeed, out of the 24 average delta scores obtained in this subgroup (table 4a, row  
403 2), 15 were negligible (<5 percentage points). Out of the remaining 9 data points, which  
404 are all negative, only two show an average decrease large enough to be considered  
405 clinically meaningful (one instance in the MPQ, and the other in the PDs).

406 The visual representation of individual delta scores (Table 4b, middle row) reveals that  
407 participants in this subgroup tended to show an uneven split between increases and  
408 decreases in scores from one visit to the next, with a larger number of individual delta  
409 scores being negative.

410

411 RTM predicts that participants with a ‘normal’ score will show little or no change from  
412 one session to the next, with an even distribution of increases and decreases cancelling  
413 each other out. This should translate in average delta scores being roughly equal to 0  
414 (Table 4a) and, visually, in a roughly even and symmetrical split between individual  
415 increases and decreases (Table 4b). As this is not the pattern of results that we observed,  
416 our results suggest that RTM alone cannot account for the overall decrease in scores seen  
417 in some outcome measures (see discussion).

418

#### 419 *Average evolution of ‘extremely small’ scores*

420 The evolution of ‘extremely low’ scorers (i.e., participants with an initial Z score <-0.5 at  
421 V1 for the V1-V2 analysis, and those with a Z score <-0.5 at V2 for the V2-V3 analysis)  
422 is shown in the third rows of Table 4a, 4b and 4c. As it can be seen in table 4a, on  
423 average, this subgroup appears to remain stable from one visit to the next on most  
424 questionnaires, with only sparse and small average increases observed on a few  
425 questionnaires. However, on roughly half of the questionnaires, these seemingly ‘stable’  
426 average deltas are a product of significant individual increases and decreases that roughly

427 cancel each other out, as shown visually in table 4b and quantified in table 4c. In the  
428 other questionnaires, the average stability does appear to stem from a widespread absence  
429 of individual fluctuations, as evidenced by the data presented in table 4b and 4c.

430 RTM predicts that extremely low scores will *increase* towards more ‘normal’ scores on  
431 the subsequent measurement. Overall, there appears to have been a slight RTM effect on  
432 a few questionnaires, although most average deltas are close to 0, suggesting that no  
433 substantial RTM was at play – or that some other effect was at play that counteracted  
434 RTM (see discussion).

435

#### 436 Analysis by questionnaire

437

##### 438 *Average evolution*

439 For the ‘extremely high’ subgroup, an overall average decrease in score was observed  
440 from one visit to the next. However, this effect was stronger in certain questionnaires:  
441 notably, the PCS, PDI, BPIi, BPIs, and MPQi all showed an average decrease larger than  
442 15 percentage points for at least one time period (Table 4b). On the other hand, other  
443 questionnaires remained comparatively more stable: the PD, CSI, and both subscales of  
444 the STAI showed no change on average from one visit to the next, for both time periods.

445

446 The ‘normal’ subgroup was more stable overall, with clinically meaningful average  
447 fluctuations (as we defined:  $\Delta < -10$ ) observed only for the MPQ and PDs, over a single  
448 time period. As with the ‘extremely high’ subgroup, the ‘normal’ subgroup was most  
449 stable on average on the CSI and both subscales of the STAI, as well as on the POQ  
450 (table 4b).

451

452 The ‘extremely low’ subgroup had the most stable scores of all, showing no clinically  
453 meaningful average change on any questionnaires, and only showing small (between 5  
454 and 10 percentage points) average changes on the PCS, PDs and BPIs, for a single time  
455 period (table 4b).

456

457 Overall, the PCS shows the largest average change for all three subgroups, followed by  
458 the PDI, PDs, both subscales of the MPQ, and both subscales of the BPI. The  
459 questionnaires with the smallest average delta were the CSI and both subscales of the  
460 STAI (table 4b).

461

#### 462 *Average absolute fluctuation*

463 As mentioned previously, it is possible for a questionnaire to have an average delta of  
464 roughly 0 from one visit to the next, seemingly suggesting that all participants remained  
465 stable over time, while in fact large individual increases and decreases in scores have  
466 been taking place, cancelling each other out. Fluctuations scores (table 4c), computed by  
467 averaging the absolute value of individual deltas within each subgroup and questionnaire,  
468 allows us to identify such cases. For example, on the PCS, the ‘normal’ subgroup had an  
469 average delta of -4 between V1 and V2. At face value, this suggests that scores remained  
470 roughly unchanged between V1 and V2 for participants in this subgroup. However, visual  
471 inspection of table 4b shows that significant increases and decreases seem to have taken  
472 place, and taking the average of the *absolute value* of these deltas allows us to quantify  
473 the average magnitude of the fluctuations – in this case, 14 percentage points (table 4c).  
474 This suggests that even for a participant whose initial score on the PCS was ‘normal’  
475 (i.e., not extreme), a fairly large change in score can be expected at the following visit. In  
476 contrast, for the ‘extremely low’ subgroup in the STAI-T – which has an equally small  
477 average delta of 3 percentage points (table 4a) – the average of the absolute value of these  
478 deltas is 4, suggesting that this questionnaire is truly stable.

479

480 For all three subgroups, the PCS shows the largest fluctuation in scores, followed closely  
481 by the MPQi, BPIs, MPQ, PDs, and BPIi. The STAI-T was the only questionnaire where  
482 participants had, on average, fluctuations smaller than 5 percentage points across visits.

483

#### 484 Healthy controls

485 HC completed the same visits and procedures as the CLBP patients but completed only  
486 the questionnaires applicable to healthy subjects: the PCS and both subscales of the  
487 STAI.

488 HC had much more stable scores overall compared to CLBP patients (Table 5). Indeed,  
489 the average evolution over time is smaller than 5 percentage points for all 3 subgroups,  
490 on all questionnaires and at both time points (with one exception at 6 percentage points)  
491 (Table 5a). Visually, apart from a notable outlier on the STAI-S (a grad student who  
492 reported having a particularly stressful day), most individual fluctuations from one visit  
493 to the next were also negligible (Table 5b). This is further supported by the average  
494 *absolute* fluctuation scores, which for the most part do not meet the threshold to be  
495 considered clinically meaningful (Table 5c). The only exception appears to be the  
496 ‘extremely high’ subgroup for the PCS, with average *absolute* fluctuations of 11 and 14  
497 percentage points from V1 to V2 and V2 to V3 (respectively). However, it is important to  
498 point out that this subgroup was comprised of only 3 participants.

499

500 Interestingly, all three questionnaires completed by the HC (the PCS, STAI-S, and STAI-  
501 T) had an average absolute fluctuation of roughly 5 percentage point across all subgroups  
502 and both time points (Table 5c). In contrast, those same questionnaires in the CLBP  
503 population had an average absolute fluctuation of 14, 9 and 4, respectively (Table 4c).

504

## 505 **Discussion**

506 The objectives of this analysis were: 1) to describe and quantify the natural trajectory of  
507 questionnaire scores over time, based on initial scores, with a subgoal of determining  
508 whether the observed fluctuations were compatible with RTM; and 2) to evaluate and  
509 compare the stability of each questionnaire over time, in 23 untreated CLBP and 25  
510 healthy controls.

511

512 Our results show that the CLBP population had relatively large variations in outcome  
513 measures over time, and that this effect varied across subgroups and across  
514 questionnaires. It bears repeating that these fluctuations were *observed in the absence of*  
515 *any experimental intervention*. Participants with high initial scores were overwhelmingly  
516 likely to show a decrease in score at the subsequent measurement, while participants with  
517 normal or extremely low scores were relatively more stable. In terms of questionnaires,  
518 the PCS showed the most variation in scores over time; both subscales of the MPQ and



519 both subscales of the BPI as well as the PDs also showed meaningful variations. The  
520 most stable questionnaire overall was the STAI-T, followed by the CSI and POQ.  
521 Healthy controls, in contrast, showed very little variability. Average deltas and average  
522 *absolute* deltas were similar - and very small - across all subgroups and questionnaires.

523

524 RTM alone cannot be responsible for all the variability observed in our CLPB sample.  
525 Indeed, RTM predicts that extreme scores are likely to be followed by less extremes  
526 (more ‘normal’) scores; as such, if RTM was the sole driving factor, extremely high  
527 scores would be followed by lower scores, and extremely low scores would similarly be  
528 followed by higher scores. However, while we did observe that extremely high scores  
529 were generally followed by lower scores, extremely low scores were *not* followed by  
530 higher scores, instead remaining relatively stable. Moreover, RTM predicts that ‘normal’  
531 scores are equally likely to show a slight increase or a slight decrease at the subsequent  
532 measurements, and that these random variations should roughly cancel out. As such, if  
533 RTM was the sole driving factor, the average delta in the ‘normal’ subgroups should be  
534 roughly 0. However, what we observed was a tendency for scores in the ‘normal’  
535 subgroup to *decrease* over time.

536

537 Together, these results suggest the presence of an effect responsible for a generalized  
538 decrease in scores (i.e., clinical improvement) over time. We hypothesize that this effect  
539 is a result of the attention and care received by the patients as part of their participation in  
540 the study, and as such propose calling this effect “Effect of Care”. Indeed, even if a  
541 participant is fully aware that they are not receiving any treatment (which therefore rules  
542 out a placebo effect, in its textbook definition<sup>29</sup>), simply having the chance to talk about  
543 their pain with understanding, thoughtful and competent-looking research staff could  
544 contribute to improve their symptoms. Additionally, the ‘seriousness’ afforded by the  
545 inclusion of brain and lumbar MRI – a notably well-regarded and imposing modality –  
546 likely further increased the potency of Effect of Care in our study.

547

548 Similarities and differences with Test-Retest

549 At first glance, this study presents superficial similarities with the well-known test-retest,  
550 such that readers might question the novelty and relevance of our findings. However, our  
551 analysis presents a number of significant differences with test-retest, and as such can  
552 offer novel and clinically relevant findings.

553

554 First, test-retest is often included as part of a validation study for a single questionnaire.  
555 Therefore, despite efforts to standardize these studies and improve generalizability, the  
556 varying study design and populations can make it difficult to compare different  
557 questionnaires. In the present study, a single patient population completed a large array of  
558 questionnaires within a fixed time frame. As such, we can easily isolate and compare the  
559 variability attributable specifically to the questionnaires.

560 Moreover, while test-retest studies generate a single overall score for a questionnaire, we  
561 conducted an analysis by subgroup. This allowed us to isolate and quantify differing  
562 degrees of variability *within* a questionnaire, and to highlight directional trends  
563 depending on the initial score, providing more nuanced and precise results than a single  
564 overall score.

565

#### 566 Biases and limitations

567 The most important limitation in this study is obviously the small sample size, especially  
568 as we further divided our sample into three subgroups. These subgroups were also of  
569 varying sizes, a result of our method of obtaining these subgroups, based on z-scores.  
570 However, having three assessment time points allowed us to conduct two separate  
571 analyses (V1 to V2, and V2 to V3) which showed similar results. Furthermore, the  
572 objective of this study was not to precisely quantify specific effects, but rather to explore  
573 a data set and identify general trends and effects. Finally, the fact that a similar pattern  
574 was found regardless of the classification threshold used (see supplementary materials)  
575 lends further credibility to our findings.

576

577 Another potentially objectionable point was the decision to use an arbitrary threshold for  
578 fluctuations that are considered ‘noise’ ( $\leq 5/100$ ) vs ‘clinically meaningful’ ( $> 10/100$ ), as  
579 opposed to using the established Minimal Detectable Change (MDC) or Minimal

580 Clinically Important Difference (MCID) of each instrument. Standardized thresholds  
581 were chosen to facilitate comparisons between questionnaires, which would otherwise  
582 have been counterintuitive at best. This decision was again consistent with our objectives,  
583 which were to identify overall trends and not to quantify phenomena with a high degree  
584 of precision (which would have been impossible given our small sample size).

585

### 586 Relevance for clinical trials

587 It is difficult to determine with certainty whether the variation in scores observed in this  
588 study are a manifestation of RTM or the result of Effect of Care or some other effect.  
589 However, being able to quantify this variability for specific questionnaires and specific  
590 subgroups has important clinical implications. Indeed, like MDC (which, unlike our, it  
591 will allow future researchers conducting clinical trials to compare their observed  
592 variations against our results, in order to isolate and better estimate the ‘true’ effect of  
593 their intervention). For example, a researcher might be thrilled to see a reduction of 10  
594 points on the PCS following an experimental treatment. However, as shown in our study,  
595 such a decrease can easily be observed in the absence of any treatment.

596

597

### 598 **Conclusion**

599

600 Ours results showed that CLBP patients with more severe symptoms at baseline will tend  
601 to show improvement at the subsequent measurement, even in the absence of an  
602 intervention – which could lead researchers to overestimate the effect of their  
603 intervention. However, patients with less severe symptoms at baseline do not show the  
604 corresponding exacerbation predicted by RTM, and patients with more average  
605 symptoms at baseline also tend to show an improvement at the subsequent measurement.  
606 These results suggest the presence of an Effect of Care, wherein patients generally show  
607 an improvement in symptoms simply by being part of a study.

608

609 Our results also provide a preliminary quantification of the variability in scores observed  
610 over time, in the absence of an intervention, in a CLBP population. This variability

611 depends on initial score and is different across questionnaires. Our results can therefore  
612 be used to guide interpretation of results obtained in clinical trials.

613

614

615

616

617

618

619

620

621

622

623

## 624 **References**

625

626 1. Kuehn B. Chronic Pain Prevalence. *JAMA*. 2018;320(16):1632.  
627 doi:10.1001/jama.2018.16009

628 2. Yong RJ, Mullins PM, Bhattacharyya N. Prevalence of chronic pain among adults in  
629 the United States. *Pain*. 2022;163(2):e328-e332.  
630 doi:10.1097/j.pain.0000000000002291

631 3. Moulin DE, Clark AJ, Speechley M, Morley-Forster PK. Chronic Pain in Canada -  
632 Prevalence, Treatment, Impact and the Role of Opioid Analgesia. *Pain Res Manag*.  
633 2002;7(4):179-184. doi:10.1155/2002/323085

634 4. Haller H, Lauche R, Sundberg T, Dobos G, Cramer H. Craniosacral therapy for  
635 chronic pain: a systematic review and meta-analysis of randomized controlled trials.  
636 *BMC Musculoskelet Disord*. 2020;21(1):1. doi:10.1186/s12891-019-3017-y

637 5. Aviram J, Samuelly-Leichtag G. Efficacy of Cannabis-Based Medicines for Pain  
638 Management: A Systematic Review and Meta-Analysis of Randomized Controlled  
639 Trials. *Pain Physician*. 2017;20(6):E755-E796.

640 6. O’Keeffe M, O’Sullivan P, Purtill H, Bargary N, O’Sullivan K. Cognitive functional  
641 therapy compared with a group-based exercise and education intervention for chronic  
642 low back pain: a multicentre randomised controlled trial (RCT). *Br J Sports Med*.  
643 2020;54(13):782-789. doi:10.1136/bjsports-2019-100780

- 644 7. Ashar YK, Gordon A, Schubiner H, et al. Effect of Pain Reprocessing Therapy vs  
645 Placebo and Usual Care for Patients With Chronic Back Pain: A Randomized Clinical  
646 Trial. *JAMA Psychiatry*. 2022;79(1):13-23. doi:10.1001/jamapsychiatry.2021.2669
- 647 8. Haack M, Simpson N, Sethna N, Kaur S, Mullington J. Sleep deficiency and chronic  
648 pain: potential underlying mechanisms and clinical implications.  
649 *Neuropsychopharmacol Off Publ Am Coll Neuropsychopharmacol*. 2020;45(1):205-  
650 216. doi:10.1038/s41386-019-0439-z
- 651 9. McWilliams LA, Cox BJ, Enns MW. Mood and anxiety disorders associated with  
652 chronic pain: an examination in a nationally representative sample. *Pain*. 2003;106(1-  
653 2):127-133. doi:10.1016/s0304-3959(03)00301-4
- 654 10. Henderson LA, Di Pietro F, Youssef AM, et al. Effect of Expectation on Pain  
655 Processing: A Psychophysics and Functional MRI Analysis. *Front Neurosci*.  
656 2020;14:6. doi:10.3389/fnins.2020.00006
- 657 11. Baird A, Sheffield D. The Relationship between Pain Beliefs and Physical and  
658 Mental Health Outcome Measures in Chronic Low Back Pain: Direct and Indirect  
659 Effects. *Healthcare*. 2016;4(3):58. doi:10.3390/healthcare4030058
- 660 12. Wideman TH, Edwards RR, Walton DM, Martel MO, Hudon A, Seminowicz DA.  
661 The Multimodal Assessment Model of Pain: A Novel Framework for Further  
662 Integrating the Subjective Pain Experience Within Research and Practice. *Clin J Pain*.  
663 2019;35(3):212-221. doi:10.1097/AJP.0000000000000670
- 664 13. Sullivan MJL, Bishop SR, Pivik J. The Pain Catastrophizing Scale: Development  
665 and validation. *Psychol Assess*. 1995;7:524-532. doi:10.1037/1040-3590.7.4.524
- 666 14. Cleeland CS, Ryan KM. Pain assessment: global use of the Brief Pain Inventory.  
667 *Ann Acad Med Singapore*. 1994;23(2):129-138.
- 668 15. Kamerman PR, Vollert J. Greater baseline pain inclusion criteria in clinical trials  
669 increase regression to the mean effect: a modelling study. *Pain*. 2022;163(6):e748-  
670 e758. doi:10.1097/j.pain.0000000000002468
- 671 16. Davis CE. The effect of regression to the mean in epidemiologic and clinical  
672 studies. *Am J Epidemiol*. 1976;104(5):493-498.  
673 doi:10.1093/oxfordjournals.aje.a112321
- 674 17. French D, Noël M, Vigneau F, French J, Cyr C, Evans R. L'Échelle de  
675 dramatisation face à la douleur PCS-CF: Adaptation canadienne en langue française de  
676 l'échelle Pain Catastrophizing Scale. / PCS-CF: A French-language, French-Canadian  
677 adaptation of the Pain Catastrophizing Scale. *Can J Behav Sci Can Sci Comport*.  
678 2005;37:181-192. doi:10.1037/h0087255

- 679 18. Anagnostis C, Gatchel RJ, Mayer TG. The Pain Disability Questionnaire: A New  
680 Psychometrically Sound Measure for Chronic Musculoskeletal Disorders. *Spine*.  
681 2004;29(20):2290-2302. doi:10.1097/01.brs.0000142221.88111.0f
- 682 19. Gauthier N, Thibault P, Adams H, Sullivan MJ. Validation of a French-Canadian  
683 version of the Pain Disability Index. *Pain Res Manag J Can Pain Soc*.  
684 2008;13(4):327-333.
- 685 20. Poundja J, Fikretoglu D, Guay S, Brunet A. Validation of the French version of  
686 the brief pain inventory in Canadian veterans suffering from traumatic stress. *J Pain*  
687 *Symptom Manage*. 2007;33(6):720-726. doi:10.1016/j.jpainsymman.2006.09.031
- 688 21. Freynhagen R, Baron R, Gockel U, Tölle TR. pain DETECT: a new screening  
689 questionnaire to identify neuropathic components in patients with back pain. *Curr Med*  
690 *Res Opin*. 2006;22(10):1911-1920. doi:10.1185/030079906x132488
- 691 22. Clark ME, Gironde RJ, Young RW. Development and validation of the Pain  
692 Outcomes Questionnaire-VA. *J Rehabil Res Dev*. 2003;40(5):381-395.  
693 doi:10.1682/jrrd.2003.09.0381
- 694 23. Spielberger C, Gorsuch R, Lushene R, Vagg P, Jacobs G. *Manual for the State-*  
695 *Trait Anxiety Inventory (Form Y1 – Y2)*. Vol IV.; 1983.
- 696 24. Gauthier J, Bouchard S. A French-Canadian adaptation of the revised version of  
697 Spielberger's State-Trait Anxiety Inventory. *Can J Behav Sci Can Sci Comput*.  
698 1993;25:559-578. doi:10.1037/h0078881
- 699 25. Melzack R. The short-form McGill Pain Questionnaire. *Pain*. 1987;30(2):191-  
700 197. doi:10.1016/0304-3959(87)91074-8
- 701 26. Boureau F, Luu M, Doubrère JF. Comparative study of the validity of four French  
702 McGill Pain Questionnaire (MPQ) versions. *Pain*. 1992;50(1):59-65.  
703 doi:10.1016/0304-3959(92)90112-O
- 704 27. Mayer TG, Neblett R, Cohen H, et al. The development and psychometric  
705 validation of the central sensitization inventory. *Pain Pract Off J World Inst Pain*.  
706 2012;12(4):276-285. doi:10.1111/j.1533-2500.2011.00493.x
- 707 28. Pitance L, Elise P, Lannoy B, et al. Cross cultural adaptation, reliability and  
708 validity of the French version of the central sensitization inventory. *Man Ther*.  
709 2016;25:e83-e84. doi:10.1016/j.math.2016.05.139
- 710 29. Turner JA, Deyo RA, Loeser JD, Von Korff M, Fordyce WE. The importance of  
711 placebo effects in pain treatment and research. *JAMA*. 1994;271(20):1609-1614.

712  
713  
714

715

716

717

718

**Table 1: Questionnaires completed by CLBP patients and HC**

Questionnaires	Subscales	Description	Items	Scale	Total score	French validated version
Pain Catastrophizing Scale (PCS) <sup>13</sup>	-	Degree of catastrophic thoughts (helplessness, magnification, and rumination)	13	5-point Likert-scale (“not at all” to “all the time”)	0 - 52	Yes <sup>17</sup>
Pain Disability Index (PDI) <sup>18</sup>	-	Ability to perform daily activities (home, social, recreational, occupational, sexual, self-care and life support activities)	7	Numerical rating scale (NRS) (0 = no disability to 10 = worst disability)	0 - 70	Yes <sup>19</sup>
Brief Pain Inventory (BPI) <i>short form</i> <sup>14</sup>	Pain severity (BPIs)	Intensity of pain (current, average, least and worst pain in the last 24h)	4	NRS (0 = no pain; 10= worst pain imaginable)	0 - 40	Yes <sup>20</sup>
	Pain interference (BPIi)	Interference of pain with daily activities (sleeping, walking, mood, etc.)	7	NRS (0 = no interference; 10=complete interference)	0 - 70	Yes <sup>20</sup>
PainDETECT (PD) <sup>21</sup>	PD	Evaluates the presence of neuropathic pain components in patients with back-pain, such as burning sensation, electric shocks etc.	9	7 items rated on a 6-point Likert Scale (0=not at all, 5=very strongly), 1 item based on pain behavior pattern score (-1, 0 or 1) & 1 item based on a radiation score (0 or 2)	0 - 38	Yes <sup>21</sup>
	Pain Detect Severity (PDs)	Intensity of pain (current, average in the past 4 weeks, worst in the past 4 weeks)	3	NRS (0 = no pain; 10= worst pain imaginable)	0 - 10	Yes <sup>21</sup>
Pain Outcomes Questionnaire (POQ) <sup>22</sup>	-	Global function (eg. mobility, vitality, affect, daily activities, pain, etc.)	19	NRS (0 = less symptoms to 10= more severe symptoms)	0 - 190	No (In-house translation)
State-Trait Anxiety Inventory (STAI-S/T) <sup>23</sup>	State anxiety (STAI-T)	Current anxiety	20	4-point Likert-scale (“not at all” to “all the time”)	20 - 80	Yes <sup>24</sup>
	State anxiety (STAI-S)	General anxiety	20	4-point Likert-scale (“not at all” to “all the time”)	20 - 80	Yes <sup>24</sup>
McGill Pain Questionnaire (MPQ) <i>short form</i> <sup>25</sup>	MPQ	Sensory-affective components of pain	15	4-point Likert scale ("no pain"; to "severe pain")	0 - 45	Yes <sup>26</sup>
	MPQ intensity (MPQi)	Pain intensity (previous week)	1	100-point visual analogue scale (left anchor: “no pain”; right anchor: “worst possible pain”)	0 - 100	Yes <sup>25</sup>
Central Sensitization Inventory (CSI) <i>short form</i> <sup>27</sup>	-	Symptoms of central sensitization	25	5-point Likert-scale (0 = “never” ; 4 = “always”)	0 - 100	Yes <sup>28</sup>

CLBP patients completed eight questionnaires: 1) Pain Catastrophizing Scale (PCS), 2) Pain Disability Index (PDI), 3) Brief Pain Inventory (BPI); 4) Pain DETECT; 5) Pain Outcomes Questionnaire (POQ), 6) State-Trait Anxiety Inventory (STAI/S-T), 7) McGill Pain Questionnaire (MPQ), and 8) Central Sensitization Inventory (short form) (CSI). HC completed only the PCS and the STAI/S-T



**Table 2: Sociodemographic characteristics of the sample.**

	CLBP (n=23)	HC (n=25)
<b>Biological sex</b>		
Women	11	15
Men	12	10
<b>Age (average ± sd)</b>	44 ± 15	40 ± 14
<b>Ethnicity</b>		
Caucasian	16	24
Asiatic	1	1
Hispanic	4	-
African	1	-
Arabic	1	-
<b>Education level</b>		
Primary school	-	-
High school	-	-
Apprenticeship	5	2
College	4	6
University	14	17
<b>Annual income</b>		
less than 20K	2	5
20K - 35K	5	2
35K - 50K	5	1
50K - 65K	5	6
65K - 80K	2	5
80K - 100K	3	4
100K and more	1	2
<b>Pain duration</b>		
4 month - 5 months	0	NA
6 months - 12 month	5	NA
1 - 4 years	7	NA
5 years and more	11	NA

**Table 3: Average scores for each questionnaire during the 3 visits (V1, V2, and V3), for the CLBP sample and the HC sample**

A)

<b>CLBP</b>	<b>PCS (0-52)</b>	<b>MPQi (0-100)</b>	<b>BPIs (0-40)</b>	<b>MPQ (0-45)</b>	<b>PDs (0-10)</b>	<b>BPIi (0-70)</b>	<b>PDI (0-70)</b>	<b>STAI/S (20-80)</b>	<b>PD (0-38)</b>	<b>POQ (0-190)</b>	<b>CSI (0-100)</b>	<b>STAI/T (20-80)</b>
<b>V1</b>	19 ± 12	56 ± 16	18 ± 5	14 ± 7	5 ± 1	18 ± 10	16 ± 9	33 ± 8	7 ± 4	44 ± 14	35 ± 11	35 ± 11
<b>V2</b>	16 ± 10	54 ± 20	18 ± 6	13 ± 7	5 ± 1	16 ± 10	14 ± 9	33 ± 9	9 ± 4	43 ± 18	33 ± 13	35 ± 12
<b>V3</b>	12 ± 11	46 ± 23	15 ± 7	10 ± 8	5 ± 2	9 ± 7	9 ± 7	32 ± 10	9 ± 5	37 ± 19	30 ± 14	35 ± 12

B)

<b>HC</b>	<b>PCS (0-52)</b>	<b>STAI/S (20-80)</b>	<b>STAI/T (20-80)</b>
<b>V1</b>	6 ± 8	26 ± 6	29 ± 9
<b>V2</b>	6 ± 7	26 ± 6	29 ± 9
<b>V3</b>	5 ± 7	26 ± 8	29 ± 10

Average scores at V1, V2 and V3 are presented in table 3a (CLPB) and 3b (HC). Scores are reported as average ± standard deviation. The theoretical min and max scores for each questionnaire are reported in the title row.

Pain Catastrophizing Scale (PCS), Pain Disability Index (PDI), Brief Pain Inventory – severity (BPIs); Brief Pain Inventory – interference (BPIi); Pain DETECT (PD); Pain DETECT severity (PDs); Pain Outcomes Questionnaire (POQ); State-Trait Anxiety Inventory (STAI/S-T); McGill Pain Questionnaire (MPQ); McGill Pain Questionnaire intensity (MPQi) and Central Sensitization Inventory (CSI) (CSI).

**Table 4: Evolution over time based on initial score – CLBP sample**

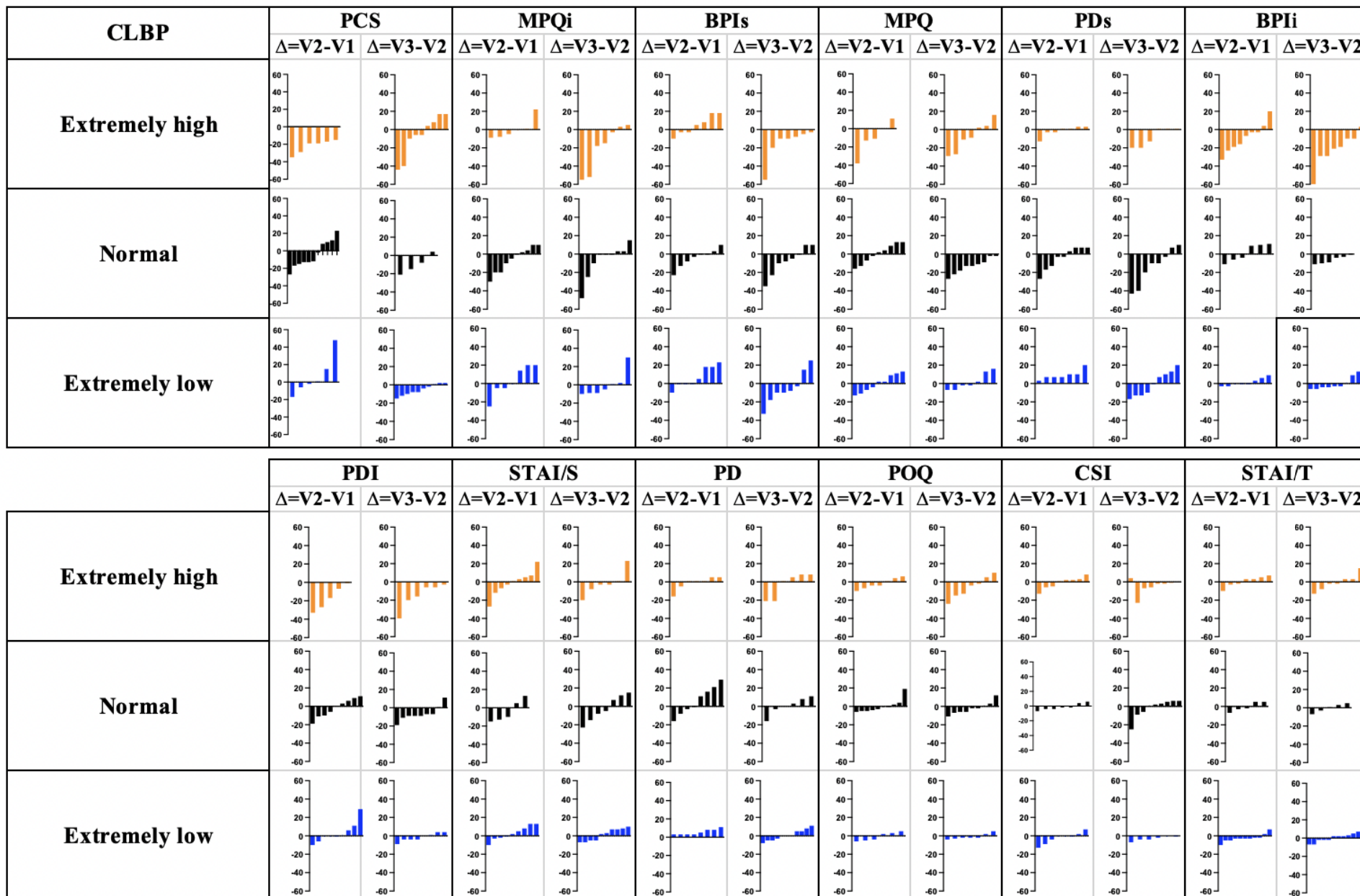
A)

CLBP	PCS (0-52)		MPQi (0-100)		BPIs (0-40)		MPQ (0-45)		PDs (0-10)		BPIi (0-70)	
	$\Delta=V2-V1$	$\Delta=V3-V2$	$\Delta=V2-V1$	$\Delta=V3-V2$	$\Delta=V2-V1$	$\Delta=V3-V2$	$\Delta=V2-V1$	$\Delta=V3-V2$	$\Delta=V2-V1$	$\Delta=V3-V2$	$\Delta=V2-V1$	$\Delta=V3-V2$
<b>Extremely high</b>	<b>-22 (-12)</b>	<b>-7 (-3)</b>	0	<b>-19 (-19)</b>	5	<b>-16 (-6)</b>	-10 (-5)	-8 (-3)	-2	-9 (-1)	-9 (-6)	<b>-22 (-15)</b>
<b>Normal</b>	-4	-10 (-5)	-6 (-6)	-7 (-7)	-4	-8 (-3)	0	<b>-13 (-6)</b>	-4	<b>-14 (-1)</b>	1	-6 (-4)
<b>Extremely low</b>	6 (3)	-5	3	0	7 (3)	-5	0	2	9 (1)	0	1	0

	PDI (0-70)		STAI/S (20-80)		PD (0-38)		POQ (0-190)		CSI (0-100)		STAI/T (20-80)	
	$\Delta=V2-V1$	$\Delta=V3-V2$	$\Delta=V2-V1$	$\Delta=V3-V2$	$\Delta=V2-V1$	$\Delta=V3-V2$	$\Delta=V2-V1$	$\Delta=V3-V2$	$\Delta=V2-V1$	$\Delta=V3-V2$	$\Delta=V2-V1$	$\Delta=V3-V2$
<b>Extremely high</b>	<b>-17 (-12)</b>	<b>-15 (-11)</b>	-1	-2	-2	-4	-2	-6 (-12)	-1	-5	0	0
<b>Normal</b>	-2	-7 (-5)	-4	-3	6 (2)	0	0	-2	-1	-2	0	0
<b>Extremely low</b>	3	-1	3	1	5	1	-1	-1	-2	-3	-3	0

B)



medRxiv preprint doi: <https://doi.org/10.1101/2023.05.05.23289575>; this version posted May 5, 2023. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY-ND 4.0 International license](#).

CLBP	PCS		MPQi		BPIs		MPQ		PDs		BPIi	
	$\Delta=V2-V1$	$\Delta=V3-V2$	$\Delta=V2-V1$	$\Delta=V3-V2$	$\Delta=V2-V1$	$\Delta=V3-V2$	$\Delta=V2-V1$	$\Delta=V3-V2$	$\Delta=V2-V1$	$\Delta=V3-V2$	$\Delta=V2-V1$	$\Delta=V3-V2$
<b>Extremely high</b>	22	17	8	22	9	16	15	14	4	9	14	23
<b>Normal</b>	14	12	11	12	7	13	9	13	10	18	9	6
<b>Extremely low</b>	15	6	13	9	9	15	8	7	9	11	3	5
<b>Total (average of all participants)</b>	16	11	11	14	8	14	10	11	8	13	9	12
<b>Total (average of both deltas)</b>	14		12		11		11		10		10	

	PDI		STAI/S		PD		POQ		CSI		STAI/T	
	$\Delta=V2-V1$	$\Delta=V3-V2$	$\Delta=V2-V1$	$\Delta=V3-V2$	$\Delta=V2-V1$	$\Delta=V3-V2$	$\Delta=V2-V1$	$\Delta=V3-V2$	$\Delta=V2-V1$	$\Delta=V3-V2$	$\Delta=V2-V1$	$\Delta=V3-V2$
<b>Extremely high</b>	17	15	9	10	5	11	5	10	5	6	5	7
<b>Normal</b>	8	9	11	12	13	7	5	5	4	7	4	4
<b>Extremely low</b>	7	4	6	6	5	5	4	3	5	3	4	3
<b>Total (average of all participants)</b>	10	9	9	9	8	7	5	6	5	5	4	4
<b>Total (average of both deltas)</b>	9		9		7		5		5		4	

Table 4 shows the evolution over time of CLBP patients. For each questionnaire, the evolution from V1 to V2 and from V2 to V3 is presented; subgroup attribution (extremely high; normal; extremely low) was based on initial scores at V1 and at V2 respectively. Deltas were calculated as V2-V1 (or V3-V2), such that a negative delta represents a decrease in score (i.e., an improvement). Table 4a shows the average delta for each subgroup; table 4b shows individual deltas within each subgroup; and table 4c shows the average of the *absolute* deltas in each subgroup.

Scores at V1, V2 and V3 were reported on 100 to facilitate comparisons between questionnaires. In table 4a, clinically meaningful deltas (>10 percentage points) are highlighted in bold, and non-trivial deltas (>5 percentage points) are accompanied by raw scores in brackets, with the min-max values of each questionnaire provided in the title row.

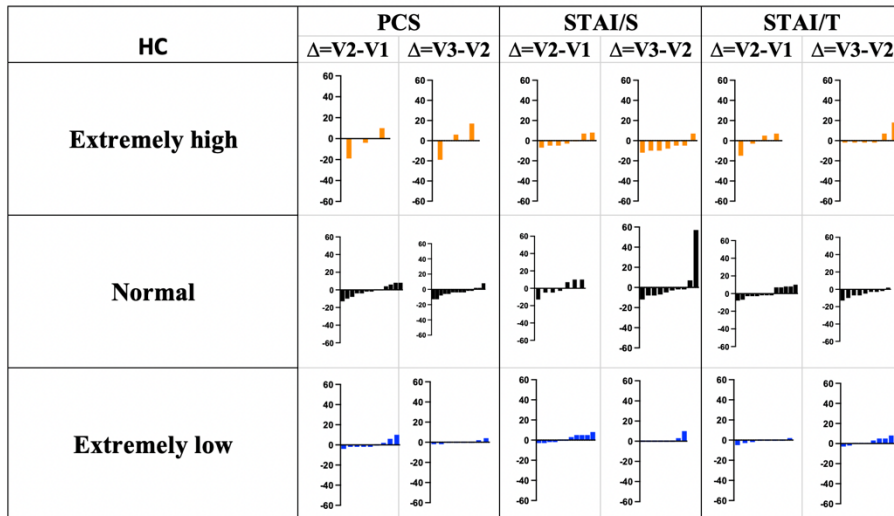
Pain Catastrophizing Scale (PCS), Pain Disability Index (PDI), Brief Pain Inventory – severity (BPIs); Brief Pain Inventory – interference (BPIi); Pain DETECT (PD); Pain DETECT severity (PDs); Pain Outcomes Questionnaire (POQ); State-Trait Anxiety Inventory (STAI/S-T); McGill Pain Questionnaire (MPQ); McGill Pain Questionnaire intensity (MPQi) and Central Sensitization Inventory (CSI) (CSI).

Table 5: Evolution over time based on initial score – HC sample

A)

HC	PCS (0-52)		STAI/S (20-80)		STAI/T (20-80)	
	$\Delta=V2-V1$	$\Delta=V3-V2$	$\Delta=V2-V1$	$\Delta=V3-V2$	$\Delta=V2-V1$	$\Delta=V3-V2$
Extremely high	-4	1	-1	-6 (-4)	-2	3
Normal	-1	-4	0	2	1	-5
Extremely low	1	0	2	2	-1	3

B)



C)

HC	PCS		STAI/S		STAI/T	
	$\Delta=V2-V1$	$\Delta=V3-V2$	$\Delta=V2-V1$	$\Delta=V3-V2$	$\Delta=V2-V1$	$\Delta=V3-V2$
<b>Extremely high</b>	11	14	5	8	8	5
<b>Normal</b>	5	5	8	11	5	6
<b>Extremely low</b>	3	1	3	2	1	4
<b>Total (average of all participants)</b>	5	5	5	7	4	5
<b>Total (average of both deltas)</b>	5		6		5	

Table 5 shows the evolution over time of HC. For each questionnaire, the evolution from V1 to V2 and from V2 to V3 is presented; subgroup attribution (extremely high; normal; extremely low) was based on initial scores at V1 and at V2 respectively. Deltas were calculated as V2-V1 (or V3-V2), such that a negative delta represents a decrease in score (i.e., an improvement). Table 5a shows the average delta for each subgroup; table 5b shows individual deltas within each subgroup; and table 5c shows the average of the *absolute* deltas in each subgroup.

Scores at V1, V2 and V3 were reported on 100 to facilitate comparisons between questionnaires. In table 5a, clinically meaningful deltas (>10 percentage points) are highlighted in bold, and non-trivial deltas (>5 percentage points) are accompanied by raw scores in brackets, with the min-max values of each questionnaire provided in the title row.

Pain Catastrophizing Scale (PCS), and State-Trait Anxiety Inventory (STAI/S-T)