

Physical fitness in third grade of primary school: A mixed model analysis of 108,295 children and 515 schools

Thea Fühner

University of Potsdam

Urs Granacher

University of Potsdam

Kathleen Golle

University of Potsdam

Reinhold Kliegl (✉ reinhold.kliegl@uni-potsdam.de)

University of Potsdam

Research Article

Keywords: physical fitness, children, primary school

Posted Date: April 1st, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-357876/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Version of Record: A version of this preprint was published at Scientific Reports on September 2nd, 2021.
See the published version at <https://doi.org/10.1038/s41598-021-97000-4>.

Physical fitness in third grade of primary school: A mixed model analysis of 108,295 children and 515 schools

Thea Fühner

University of Potsdam

Urs Granacher

University of Potsdam

Kathleen Golle

University of Potsdam

Reinhold Kliegl (✉ reinhold.kliegl@uni-potsdam.de)

University of Potsdam

Research Article

Keywords: physical fitness, children, primary school

DOI: <https://doi.org/10.21203/rs.3.rs-357876/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

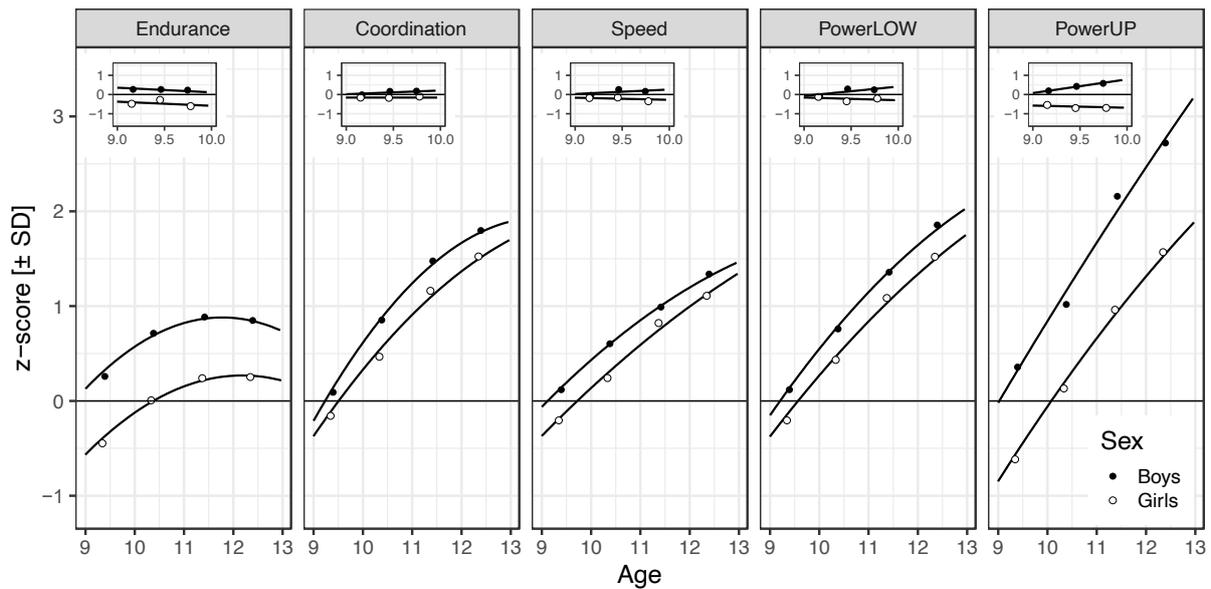
[Read Full License](#)

1 **Abstract**

2 **Children’s physical fitness development and related moderating effects of age and sex are well**
3 **documented, especially boys’ and girls’ divergence during puberty. The situation might be different**
4 **during prepuberty. As girls mature approximately two years earlier than boys, we tested a possible**
5 **convergence of performance with five tests representing four components of physical fitness in a**
6 **large sample of 108,295 eight-year old third-graders. Within this single prepubertal year of life and**
7 **irrespective of the test, performance increased linearly with chronological age, and boys**
8 **outperformed girls to a larger extent in tests requiring muscle mass for successful performance.**
9 **Tests differed in the magnitude of age effects (gains), but there was no evidence for an interaction**
10 **between age and sex. Moreover, “physical fitness” of schools correlated at $r = 0.48$ with their age**
11 **effect which might imply that “fit schools” promote larger gains; expected secular trends from**
12 **2011 to 2019 were replicated.**

13 **Introduction**

14 Children’s development of physical fitness as well as the effects of moderating variables such as age
15 and sex are well documented¹⁻⁷, especially boys’ and girls’ divergence during puberty starting late in
16 the twelfth and tenth year of life, respectively⁸. The situation might be different during prepuberty.
17 Girls mature visibly about two years earlier than boys, but sex hormones rise already much starting
18 at eight years of age⁹⁻¹¹. If the early rise of sex hormones relates to body composition, the question
19 arises whether there is evidence for faster development of girls than boys in this year, leading to a
20 convergence of performance in physical fitness during the transition from pre-puberty to puberty?
21 A longitudinal study on the development of different fitness components highlights sex-specific
22 performance trajectories in youth aged 9 to 12 years⁶. This study was a precursor project of the
23 present cross-sectional study and tested four components of physical fitness in four annual
24 assessments including 240 children⁶. The main panels in Figure 1 illustrate sex- and age-differential
25 development for cardiorespiratory endurance, coordination, speed, and power (assessed separately
26 for lower [powerLOW] and upper [powerUP] limbs). The five tests used for the assessment are
27 described in the figure caption.



28

29 **Figure 1.** Physical fitness curves of a longitudinal sample of 240 German boys (closed circles) and girls (open circles)
 30 followed from age 9 to 12 years for endurance = cardiorespiratory endurance (i.e., 9 min run test), coordination
 31 (i.e., running in a star like pattern), speed (i.e., 50-m linear sprint test), powerLOW = power of lower limbs (i.e.,
 32 triple hop test), and powerUP = power of upper limbs (i.e., ball push test). The insets show for each test score the
 33 regression on age for the first assessment when children were between 9.00 and 9.99 years old. Also shown are
 34 the means for groups of boys and girls binned into three age groups (i.e., 9.00-9.33; 9.34-9.66; 9.67-9.99). Error
 35 bands are 95% Cis. Data are from Golle et al.⁶.

36 Speed, powerLOW, and powerUP follow a mostly linear trajectory, whereas cardiorespiratory

37 endurance and coordination are characterized by a curvilinear development⁶. These age-related

38 differences can mostly be attributed to growth (increasing body mass and body height) and

39 maturation during childhood and adolescence¹. For instance, while increased body mass may have a

40 positive impact on ball push test performance, it may negatively influence performance in tests such

41 as the 6 or 9 min run tests which afford the continuous acceleration of the body¹². Furthermore, the

42 fact that boys significantly outperform girls in these tests⁶ is most likely due to differences in muscle

43 mass favouring boys^{1,13}. In fact, there is evidence that prepuberal boys have on average 3.7% larger

44 muscle mass compared with girls. In line with this argument, Golle et al.⁶ observed larger sex

45 differences in fitness tests demanding muscle mass (e.g., powerUP) compared with tests of motor

46 coordination (e.g., running in a starlike pattern).

47 While the longitudinal study of Golle et al.⁶ showed performance increases for all components of

48 physical fitness over a period of four years, it may come as a surprise that none of the *cross-sectional*

49 age differences and none of the interactions with sex were significant *within* the respective

50 assessment years - neither when aggregated over tests nor when tested for individual tests; the
51 insets in Figure 1 show these non-significant trends and the means for three age groups for the first
52 assessment of tests when children were between 9.00 and 9.99 years (see Supplement B for details).
53 However, the absence of evidence for cross-sectional age differences and their interactions with sex
54 within a single year of life at a prepubertal development stage must not be taken as evidence for
55 their absence¹⁴, because statistical power may not have been sufficiently large..

56 Our study included a very large representative sample of 108,295 eight-year old third graders. This
57 large sample allowed us to zoom into a cross-sectional *short-term ontogenetic* window. In this single
58 year of life, we expected to detect age differences that were not significant in the insets of Figure 1.
59 Indeed, we expected that the component differences in developmental gains (i.e., the differences in
60 slopes) will be in general agreement with those in the main panel of Figure 1 because gains should be
61 larger (favouring older children) in tests that do not demand continuous acceleration of the body as
62 with the ball push test vs the 6 min run test¹². Similarly, as in Figure 1 (both in the main panel and in
63 the insets) the test-related sex effects (i.e., the differences between the lines) will be larger, the
64 more muscle mass is required to perform a physical fitness test^{1,13} (favouring boys) and smaller for
65 tests involving motor coordination¹⁵. In other words, in general, we expected to anticipate the *long-*
66 *term ontogenetic longitudinal trends* across four years with *short-term ontogenetic cross-sectional*
67 *trends* within a single year for younger children in the third grade.

68 There is one exception to this expectation of agreement between *short-term* and *long-term*
69 *ontogenetic* trends: The only significant interaction between age and sex was reported for powerUP
70 (see right panel in Figure 1) indicating divergence between boys and girls. If physical fitness
71 components carry an early prepubertal signal in eight-year olds, then we should observe
72 convergence of scores because girls will benefit earlier than boys from the rise of sex hormones^{9,11}.

73 Cross-sectional analysis has been criticized for good reasons and in general preference is given to
74 longitudinal analysis, mostly because only the latter delivers information about intraindividual
75 growth⁶. However, we propose that, when the focus is on *short-term ontogenesis*, that is on changes

76 *within*, not *between* years of life, a cross-sectional design is probably the only option to determine
77 development-related gains because this design circumvents practice effects, a necessary
78 consequence of repeated testing of physical fitness within a year. Another problem of cross-sectional
79 designs are cohort effects (i.e., age-correlated cultural change). However, cohort effects are certainly
80 negligible for children who attend the same grade but differ in age only within the same year of life.
81 Thus, with the exception of loss of information about interindividual differences in intraindividual
82 change, a cross-sectional design is more suitable to determine how within-year developmental
83 profiles differ between *health-related* (e.g., cardiorespiratory endurance)¹⁶ and *skill-related*
84 components of physical fitness (e.g., power, speed, coordination)¹⁶.

85 Just as there are individual differences in physical fitness between children, there are also differences
86 between the over 500 schools in how much they implement programs that facilitate gains in the
87 development of physical fitness¹⁷. Explanatory hypotheses about these differences are beyond the
88 scope of this article, but our use of a linear mixed model for statistical inference affords exploratory
89 tests for their presence and adjusts test statistics for school-related sources of variance. Finally, as
90 data collection for this cross-sectional study occurred annually from 2011 to 2019, secular trends are
91 another source of variance in scores that must be taken into account. Here we expected that the
92 results will be in line with recent original research and meta-analyses and show a decline in
93 cardiorespiratory endurance and an increase in speed^{18–20}.

94 In summary, the hypotheses about age- and sex-related differences in physical fitness tests were
95 tested with a linear mixed model (LMM) that afforded the simultaneous consideration of children,
96 schools, and cohorts as random factors and the estimation of variance components and correlation
97 parameters for (a) test scores, (b) effects of contrasts between the tests, and (c) in the case of
98 schools and cohorts also effects of age-related gains and sex differences. These model parameters
99 serve primarily as measures of statistical control for the fixed effect estimates but may also yield
100 substantive insights about the dynamics of development. For example, for correlations of scores, we
101 expected the usual positive manifold between the tests, but for correlations of effects (i.e., the

102 contrasts between tests) no directed hypotheses were formulated given the usually low reliability
103 associated with difference scores. Although age and sex are between-child factors, they are also both
104 within-school and within-cohort factors, providing us with the opportunity to detect reliable variance
105 components and correlation parameters for these factors. No directed hypotheses were formulated
106 for schools. However, secular trends are well documented, and we expected to replicate them.

107 **Results**

108 ***Overview***

109 Table 1 displays statistics for fixed effects of age (linear) and sex as well as their interactions with the
110 four test contrasts for LMM $m2$. Test-specific z-transformations eliminated main effects of contrasts
111 H1 to H4 (all $z \leq 1.27$). Neither the age x sex interaction nor any of the interactions of this term with
112 the four test contrasts were significant (all $|z| \leq 1.74$; $p > 0.08$). Adding quadratic trends of age and
113 their associated interactions to the model, did not significantly contribute to goodness of fit; $\chi^2(10) =$
114 9.17 , $p = 0.52$. None of the age x sex interaction was significant when tested separately for the five
115 tests (all $|z| \leq 1.04$, $p > 0.29$; see Supplement A for details about both control LMMs).

116 ***Age-related gains***

117 Figure 2 displays both the gains in physical fitness with age and that boys scored higher than girls in
118 each of the five physical fitness tests. The parallel lines in each panel also visualize the lack of
119 significant evidence for age x sex as well as age x sex x test interactions. Counter to folklore, a very
120 large data set (i.e., 108,295 subjects and 525,126 observations) did not “automatically” render
121 everything significant! Given this statistical power, we are strongly inclined to interpret the absence
122 of evidence for interactions as evidence of absence of interactions for these five tests of physical
123 fitness¹⁴.

Table 1 Fixed-effect estimates of linear mixed model

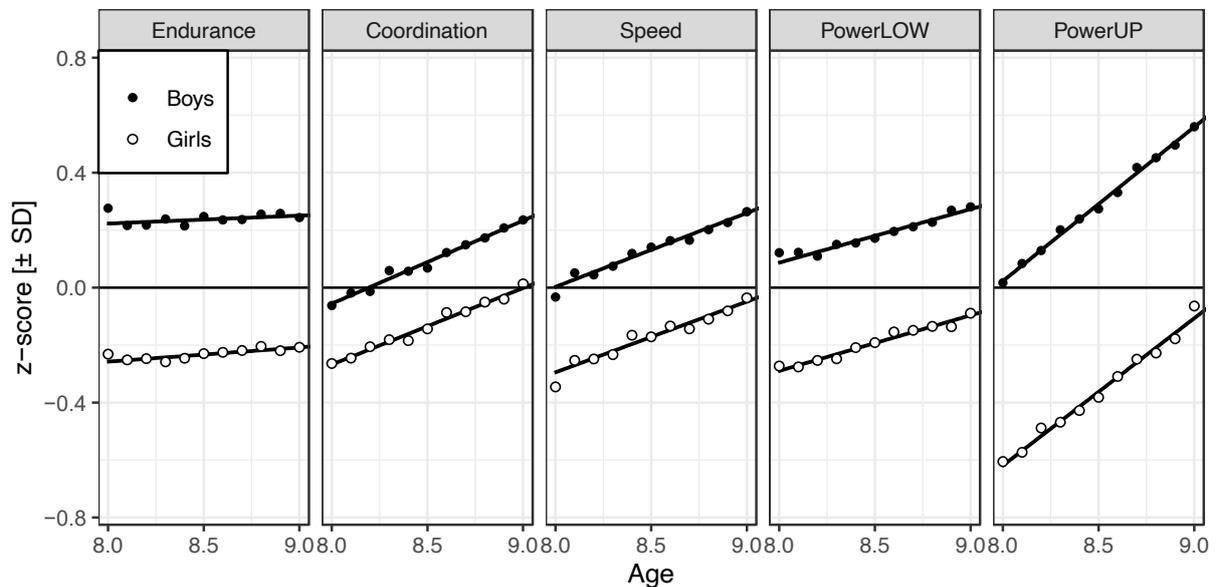
Source of variance	Fixed-effect estimates	Standard error	z-values	Pr (> z)
Main effects				
Grand mean (intercept)	-0.041	0.010	-4.00*	0.000
H1: coordination vs. endurance	0.009	0.023	0.40	0.691
H2: speed vs. coordination	-0.031	0.031	-1.01	0.310
H3: powerLOW vs. speed	0.040	0.032	1.27	0.210
H4: powerUP vs. powerLOW	-0.007	0.019	-0.34	0.731
Age (linear)	0.262	0.008	31.74*	< 0.001
Sex	0.398	0.005	86.61*	< 0.001
Age (linear) x Sex	0.002	0.013	0.14	0.887
Age (linear) x Test				
H1: coordination vs. endurance	0.207	0.011	18.64*	< 0.001
H2: speed vs. coordination	-0.064	0.011	-6.02*	< 0.001
H3: powerLOW vs. speed	-0.002	0.010	-0.24	0.808
H4: powerUP vs. powerLOW	0.298	0.011	25.68*	< 0.001
Sex x Test				
H1: coordination vs. endurance	-0.243	0.006	-38.39	< 0.001
H2: speed vs. coordination	0.077	0.006	12.71*	< 0.001
H3: powerLOW vs. speed	0.071	0.006	12.40*	< 0.001
H4: powerUP vs. powerLOW	0.287	0.007	43.38*	< 0.001
Age (linear) x Sex x Test				
H1: coordination vs. endurance	0.037	0.021	1.74	0.082
H2: speed vs. coordination	-0.012	0.020	-0.63	0.530
H3: powerLOW vs. speed	-0.019	0.019	-0.99	0.322
H4: powerUP vs. powerLOW	0.024	0.022	1.08	0.282

125 H1 to H4 = hypothesis 1 to 4, endurance = cardiorespiratory endurance (i.e., 6 min run test), coordination = star run test,
 126 speed = 20-m linear sprint test, powerLOW = power of lower limbs (i.e., standing long jump test), powerUP = power of upper
 127 limbs (i.e., ball push test), * = z-value > 3.0, linear mixed model random factors: cohorts (9), schools (515), children (108,295),
 128 observations = 525,126 (missing = 3%). For estimates of variance components and correlation parameters see Table 3.

129 Despite the small age range, the differences were large and visible for all five physical fitness tests;
 130 the overall linear trend for age was significant with $b = 0.26$, $z = 31.7$. The LMM tested the
 131 interactions of age with the four test contrasts, that is whether slopes in neighbouring panels
 132 (averaged across sex) were parallel. Three of four expected interactions were significant (see second
 133 block of Table 1): the age effect was larger for coordination than cardiorespiratory endurance (H1; b
 134 = 0.21, $z = 18.6$), larger for coordination than speed (H2; $b = -0.06$, $z = -6.0$), and larger for powerUP
 135 than powerLOW (H4; $b = 0.30$, $z = 25.7$). The difference between age slopes for powerLOW and
 136 speed (H3) was not significant ($b = -0.00$, $z = -0.2$, $p = 0.808$).

137 Counter to this profile of differences, Figure 2 which is based on observed scores suggests that speed
 138 gain is larger than powerLOW gain. Indeed, this contrast was significant for LMM $m1$, that is as long
 139 as VCs for cohort-related differences between tests were not in the LMM. The large differences in

140 secular trends are shown in Figure 4 (below); adjusting for them in the LMM yielded the reported
 141 partial effects.



142
 143 **Figure 2.** Performance differences between 8.0 and 9.0 years by sex in the five physical fitness tests presented as z-
 144 transformed data computed separately for each test. Endurance = cardiorespiratory endurance (i.e., 6 min run
 145 test), Coordination = star run test, Speed = 20-m linear sprint test, PowerLOW = power of lower limbs (i.e.,
 146 standing long jump test), PowerUP = power of upper limbs (i.e., ball push test), SD = standard deviation. Points are
 147 binned observed child means; lines are simple regression fits to the observations; 95% confidence intervals for
 148 means ≈ 0.05 are not visible.

149 **Sex-related effects**

150 The difference between lines in Figure 2 displays the expected differences between boys and girls for
 151 the performance in the five physical fitness tests; the overall sex effect was estimated with $b = 0.40$, z
 152 $= 86.6$. The third block of Table 1 lists statistics for the interactions between sex and the tests
 153 contrasts. All interaction terms were significant and in agreement with *a priori* expectations. Boys
 154 performed better than girls on cardiorespiratory endurance than coordination (H1; $b = -0.24$, $z = -$
 155 38.4), better on speed than coordination (H2; $b = 0.08$, $z = 12.7$), better on powerLOW than speed
 156 (H3; $b = 0.07$, $z = 12.4$), and better on powerUP than powerLOW (H4; $b = 0.29$, $z = 43.4$). The
 157 magnitude of these interactions is shown in Figure 2 in the differences between the parallel lines for
 158 boys and girls for neighbouring panels (averaged over age).

159 **Variance components and correlation parameters**

160 **Test scores**

161 Table 2 lists estimates of VCs and CPs for the five test scores from a re-parameterized version of
 162 LMM `m2` with the same goodness of fit and the same estimates for fixed-effects. The test-related
 163 VCs were large children (0.69 to 0.77), of medium-size for schools (0.23 to 0.36), and small for
 164 cohorts (0.03 to 0.06). VCs for the age-related gains (slopes 0.09) and the sex effect (0.05) were also
 165 small for schools. It is noteworthy that the differences between schools in the age-related gain of
 166 their children is larger than the differences between cohorts.

167 **Table 2** Variance components, correlation parameters, and zero-order correlations for test scores

Test	VC		CP \ r					
	End	Coord	Speed	PowerLOW	PowerUP	Age	Sex	
Child								
Endurance	0.69	1.00	0.37	0.41	0.42	0.23		
Coordination	0.72	0.57	1.00	0.44	0.44	0.32		
Speed	0.73	0.63	0.67	1.00	0.52	0.31		
PowerLOW	0.77	0.60	0.66	0.77	1.00	0.37		
PowerUP	0.70	0.25	0.45	0.43	0.50	1.00		
School								
Endurance	0.30	1.00	0.33	0.33	0.38	0.17	0.15	0.05
Coordination	0.36	0.36	1.00	0.32	0.37	0.25	0.14	0.03
Speed	0.30	0.34	0.33	1.00	0.46	0.29	0.10	0.06
PowerLOW	0.25	0.37	0.39	0.43	1.00	0.30	0.12	0.06
PowerUP	0.23	0.19	0.19	0.28	0.27	1.00	0.13	-0.04
Age	0.09	0.45	0.34	0.29	0.32	0.23	1.00	0.10
Sex	0.05	0.12	-0.05	0.11	0.19	-0.03	0.26	1.00
Cohort								
Endurance	0.05							
Coordination	0.03							
Speed	0.06							
PowerLOW	0.04							
PowerUP	0.03							

168 End = cardiorespiratory endurance (i.e., 6 min run test), Coord = star run test, Speed = 20-m linear sprint test, PowerLOW =
 169 power of lower limbs (i.e., standing long jump test), PowerUP = power of upper limbs (i.e., ball push test), VC = variance
 170 component, CP \ r = correlation parameter \ zero-order correlation; linear mixed model correlation parameters are shown
 171 below and corresponding pairwise zero-order correlations above the diagonal for children (top) and schools (middle).
 172 Theoretically relevant correlations are set in **bold**. VC for Residual = 0.53. VCs and CPs are based on full set of data; ZOC
 173 correlations are based on subsets of data (Child: 96,529 children from 512 schools; School: 93,661 children from 421
 174 schools).

175 In the first block of Table 2, CPs between tests scores for children are listed below the diagonal and
 176 the corresponding zero-order correlations (ZOCs) above the diagonal of the respective correlation
 177 matrix. The ZOCs are based on 96,529 children with complete test scores; they came from 512

178 different schools. In the second block, we list corresponding results for schools. For this analysis, we
179 also added the criterion that a school had to report complete data from more than 30 boys and 30
180 girls to ensure stable estimation of age and sex effects within schools. This criterion left us with 421
181 schools, 93,661 children, and 468,305 scores.

182 There were three noteworthy patterns of results. First, as expected, all child- and school-related CPs
183 and ZOCs between test scores were positive. Thus, the five tests represent a latent construct
184 “physical fitness” both for differences between children and for differences between schools. This
185 was also supported by a random-effects principal component analyses (rePCA) of the two orthogonal
186 random-effect structures. The first principal component (PC1) loadings ranged from 0.49 to 0.34 for
187 the child-related PC1 and from 0.45 to 0.29 for the school-related PC1, accounting for 65% and 38%
188 of the respective variances (see Supplement A for details).

189 Second, the tests did not correlate equally highly with each other. Most notably, the child-related CPs
190 of cardiorespiratory endurance, coordination, speed, and powerLOW correlated very highly between
191 0.57 and 0.77, but their correlations with powerUP were distinctly smaller (CPs: 0.25 to 0.50). This
192 observation holds also for the other three correlation matrices, but overall correlations were smaller
193 (see Table 2 for all CPs).

194 Again, this interpretation was supported by rePCAs (see Supplement A for details). The smallest
195 loading on child- and school-related PC1s was obtained for powerUP (0.34, 0.29). Moreover, for
196 children the loadings for the second PC2s (15%) represented the difference between
197 cardiorespiratory endurance (0.51) and powerUP (-0.85). Similarly, for schools the third PC3 (11%)
198 represented the difference between the average of cardiorespiratory endurance (0.32) and
199 coordination (0.45) and powerUP (-0.74). We took this PC2/PC3-based difference score as support
200 for the hypothesis that powerUP favoured heavier children for who strength of arms may mask
201 reduced cardiorespiratory endurance and coordination (see Discussion).

202 Third, child-related CPs (Table 2, top panel, below diagonal) were larger than child-based ZOCs (top,
203 above diagonal). This was a rather striking pattern because one might expect the opposite given that
204 ZOCs were confounded with large effects of age and sex. Conversely, CPs were larger despite
205 adjustment for sex and age differences in the fixed effects and for differences due to schools and
206 cohorts in the random-effect structure of the LMM. The reason for the result is LMM-based
207 shrinkage of conditional means of the units of the random factors in the direction of the GM. Thus,
208 the entire data set was used to “correct” unreliable observations, also called “borrowing strength” to
209 improve predictions. Due to this shrinkage correction for unreliability, CPs revealed the latent
210 relations between measures much more clearly than ZOCs.

211 ***Effects of test contrasts***

212 In the random-effect structure of LMM *m2*, estimates were returned for child-, school-, and cohort
213 related VCs for GM and the four test contrasts; VCs of age and sex were also estimated for school.
214 CPs for child and school reflect correlations between the contrasts (i.e., effect correlations). The
215 results are shown in Table 3. As in Table 2, CPs are reported below the diagonals and corresponding
216 ZOCs above the diagonals.

217 VCs for test contrasts were larger (0.48 to 0.72) for children and somewhat smaller, but still highly
218 reliable (0.29 to 0.38) for schools, especially when compared to VCs estimated for school-related age
219 (0.09) and sex (0.05) effects, and especially when compared to cohort-related effects (0.04 to 0.09).
220 CPs and ZOCs of effects are smaller than CPs and ZOCs based on test scores because, with the
221 exception of those involving GM, they are all based on difference scores.

222 There were two results of theoretical relevance. First, there was a negative CP for $r_{GM,H4}$ for children
223 (-0.31): If we invert the difference score to convert to a positive correlation, then large values of GM
224 correspond to a large difference between powerLOW and powerUP. In other words, the larger
225 powerLOW relative to powerUP, the larger is the expected value for GM. Thus, in line with the
226 special status of powerUP reported above, powerUP is better thought of as an adjustment or

227 sharpening of physical fitness as indicated by PowerLOW than as a genuine indicator of physical
 228 fitness by itself. The corresponding ZOC was only -0.13.

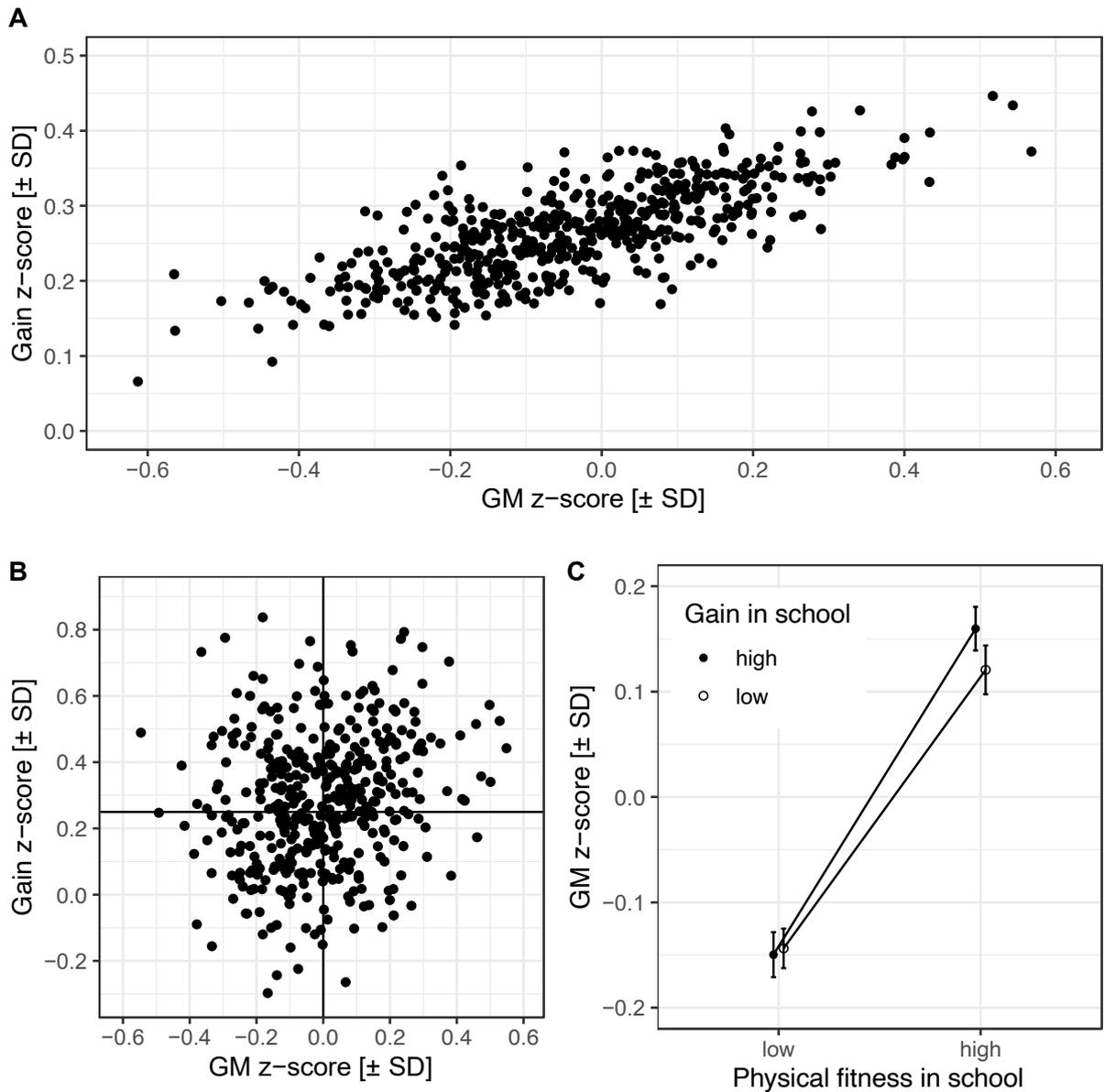
229 **Table 3** Variance components, correlation parameters, and zero-order correlations for test-related contrasts

Effects	VC				CP \ r			
	GM	H1	H2	H3	H4	Age	Sex	
Child								
Grand Mean	0.58	1.00	0.04	0.04	0.04	-0.13		
H1: coordination vs. endurance	0.64	0.13	1.00	-0.51	-0.01	0.05		
H2: speed vs. coordination	0.57	0.07	-0.51	1.00	-0.46	-0.07		
H3: powerLOW vs. speed	0.48	0.08	0.02	-0.35	1.00	-0.40		
H4: powerUP vs. powerLOW	0.72	-0.31	0.14	-0.19	-0.28	1.00		
School								
Grand Mean	0.19	1.00	0.09	-0.05	-0.09	-0.20	0.19	0.05
H1: coordination vs. endurance	0.38	0.15	1.00	-0.55	0.00	0.05	0.01	-0.01
H2: speed vs. coordination	0.38	-0.13	-0.59	1.00	-0.50	-0.02	-0.05	0.02
H3: powerLOW vs. speed	0.29	-0.11	0.02	-0.50	1.00	-0.39	0.00	-0.00
H4: powerUP vs. powerLOW	0.29	-0.19	-0.05	0.06	-0.39	1.00	0.00	-0.09
Age (linear)	0.09	0.48	-0.03	-0.10	-0.01	-0.10	1.00	0.10
Sex	0.05	0.09	-0.14	0.14	0.05	-0.19	0.26	1.00
Cohort								
Grand Mean	0.01							
H1: coordination vs. endurance	0.05							
H2: speed vs. coordination	0.08							
H3: powerLOW vs. speed	0.09							
H4: powerUP vs. powerLOW	0.04							

230 H1 to H4 = hypothesis 1 to 4, endurance = cardiorespiratory endurance (i.e., 6 min run test), coordination = star run test,
 231 speed = 20-m linear sprint test, powerLOW= power of lower limbs (i.e., standing long jump test), powerUP = power of upper
 232 limbs (i.e., ball push test), VC = sqrt (variance component), CP \ r = correlation parameter \ zero-order correlation; linear
 233 mixed model correlation parameters are shown below and corresponding pairwise zero-order correlations above the
 234 diagonal for children (top) and schools (middle). Theoretically relevant correlations are set in **bold**. VC for Residual = 0.54.
 235 VCs and CPs are based on full set of data; ZOC correlations are based on subsets of data (Child: 96,529 children from 512
 236 schools; School: 93,661 children from 421 schools).

237 Second, we note three large negative CPs (r H1.H2, r H2.H3, and r H3.H4). However, these
 238 correlations are ambiguous because the contrasts had a test in common (i.e., coordination is part of
 239 H1 and H2; speed is part of H2 and H3; powerLOW is part of H3 and H4).

240 Third, the largest effect CP was observed between age and GM for schools (+0.48) suggesting that
 241 “fitter” schools promote more developmental change across the school year. Figure 3A displays a
 242 visualization of the CP using a scatterplot of the conditional means of age-related gain over
 243 conditional means of GMs of physical fitness for the 515 schools.



244
245
246
247
248
249
250

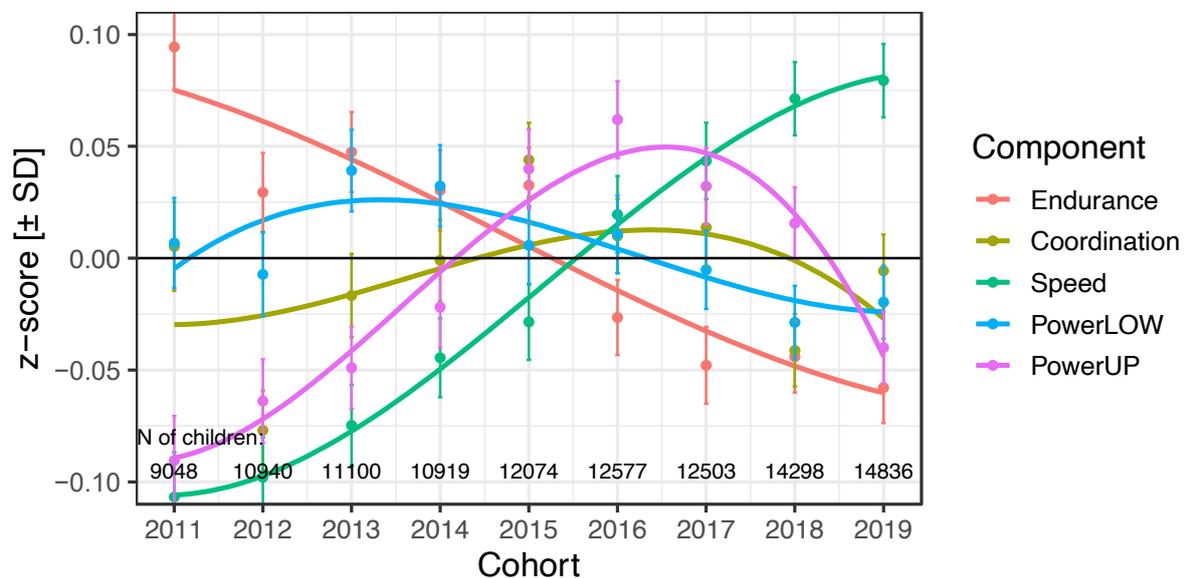
Figure 3. (a) Visualization of correlation parameter between Grand Mean (GM) and yearly gains (age effect) of 515 schools using conditional means resulting from shrinkage correction of observed data with LMM parameters. (b) Scatterplot of observed GMs and yearly gains (within-school age slopes) for 421 schools reporting data from more than 30 boys and 30 girls. (c) After sorting schools into groups of high vs. low physical fitness (split at 0) and high vs. small gain-rate (split at the gain of +0.25), the interaction corresponding to the CP becomes visible; error bars are 95% CIs. SD = standard deviation.

251 The scatterplot is not identical with the CP, in fact the correlation is 0.75, because conditional means
252 are “predictions” of the school age effect and GM using the school data and all model parameters to
253 correct for unreliability in the scores. Indeed, as shown in Figure 3B, there is no evidence for this
254 relation in the uncorrected scatterplot corresponding to the ZOC of 0.19. A significant CP
255 corresponds to a simple interaction in the data and this interaction can be visualized by sorting
256 schools into those of high and low physical fitness (split at z-score = 0) and those with a high and

257 small gain for their children in the third grade (split at the age slope of +0.25; see Figure 3C). In an
 258 ANOVA of schools' observed mean physical fitness values the interaction between these two *post-*
 259 *hoc* grouping factors was significant; $F(1, 417) = 4.22$, $MSe = 0.012$, $p < 0.05$. Of course, from a
 260 correlation we cannot infer the direction of causality. "Fitter schools" (e.g., schools offering
 261 extracurricular sport-related activities) may facilitate gains in children's fitness. Alternatively, if we
 262 assume that fitter children gain more in a year, then a school's high fitness as well as the associated
 263 large gain could be the result of being attended by fitter children (e.g., due to the school's location in
 264 a high-SES region).

265 **Cohort-related variance components**

266 The small, but significant VCs related to the random factor cohort indicate that there were reliable
 267 test x cohort interactions; they are shown in Figure 4. Across the nine years from 2011 to 2019 there
 268 was a performance decline for cardiorespiratory endurance and an increase for speed. The other
 269 three components exhibit an initial increase followed by a decline of performance in recent years.



270
 271 **Figure 4.** Cohort-related change of components of physical fitness. Points are observed means with 95% CIs. Lines are third-
 272 order polynomial trends fitted to children's scores along with 95% error. Note the much smaller range of the y-
 273 axis (i.e., from -0.10 to +0.10) compared to age effects shown in Figure 2. Endurance = cardiorespiratory
 274 endurance (i.e., 6 min run test), Coordination = star run test, Speed = 20-m linear sprint test, PowerLOW = power
 275 of lower limbs (i.e., standing long jump test), PowerUP = power of upper limbs (i.e., ball push test).

276 ***Goodness-of-fit statistics and model residuals***

277 Additional details about the LMM analyses are documented in Supplement A which also contains
278 information about the control LMMs. Despite their complexity, all models converged without
279 problems and there was no evidence of estimates of parameters at their boundaries.
280 Overparameterization was observed only for the most complex LMM *m4*. Thus, with this exception,
281 the LMMs were supported by the data. Finally, we carried out residual-based diagnostics (e.g., q-q
282 plot, standardized residuals over fitted values, etc.) for the reference LMM *m2* with CPs for effects
283 (Table 1, Table 3). These tests did not reveal any problems.

284 **Discussion**

285 The aim of this study was to examine *short-term ontogenetic cross-sectional* developmental
286 differences in physical fitness for five tests tapping *health-* and *skill-related* components of physical
287 fitness in a large sample of 108,295 German eight-year old children. Even in a single prepubertal year
288 of life (1) performance increases linearly with chronological age in all physical fitness tests, (2) boys
289 outperform girls in all physical fitness tests with sex differences being larger for tests requiring
290 muscle mass and being smaller for tests requiring motor coordination, (3) the tests differ, mostly as
291 expected, in the size of age and sex effects, (4) four of the five tests represent a common construct
292 (i.e., correlate strongly positively with each other) – the exception is the ball push test that requires
293 powerUP, (5) there was no evidence for an interaction of age and sex – with each other or with the
294 test contrasts despite an abundance of statistical power, (6) “physically fit schools” apparently
295 promote more developmental gains within a year (but this is only correlational evidence) and (7)
296 diverging secular trends for cardiorespiratory endurance (negative) and speed (positive) are in
297 agreement with other research and meta-analyses.

298 Boys outperformed girls in all four physical fitness components (cardiorespiratory endurance,
299 coordination, speed, power [LOW/UP]). This is in line with other studies reporting normative values²⁻

300 ^{5,21–25}. For instance, Tambalis et al.⁵ reported that boys aged 6 to 18 years showed significantly better
301 performances for cardiorespiratory endurance (i.e., 20 m shuttle run test), powerLOW (i.e., standing
302 long jump test), and agility (i.e., 10 x 5 m agility shuttle run test) compared to girls. Furthermore, De
303 Miguel-Etayo et al.³ described that boys aged 6 to 10 years significantly outperformed girls in
304 cardiorespiratory endurance (i.e., 20 m shuttle run test), speed (i.e., 40 m sprint test), and
305 powerLOW (i.e., standing long jump test).

306 What are the reasons for the observed large test-specific sex differences? Sex-related differences in
307 body composition appear to be likely candidates to account for the observed findings. The larger sex
308 effect in powerUP compared to powerLOW can be explained by a better proportion of strength
309 relative to body mass especially in the upper limbs in boys compared to girls^{26,27}. The large difference
310 in cardiorespiratory endurance can be explained by physiological factors such as boys' better
311 mechanical efficiency and fractional utilisation of oxygen^{2,28}. Furthermore, muscle mass¹ and muscle
312 cross-sectional area¹³ favour boys especially in physical fitness tests that recruit muscle mass. Beside
313 these anthropometric factors and physiological demands, sociocultural aspects may also explain the
314 sex difference. Haywood and Getchell²⁹ reported that girls usually participate in sports that require
315 balance and flexibility (e.g., gymnastics, figure skating) compared to boys who rather participate in
316 strength-related activities.

317 The sex effect was also significantly stronger for powerLOW than for speed. PowerLOW is
318 determined much more by muscle mass where boys usually outperform girls^{2,3,5}. In contrast, speed is
319 less influenced by muscle mass than by motor coordination where sex differences are comparatively
320 small³⁰ or were not found at all¹⁵. Therefore, the sex effect in powerLOW might be larger than in
321 speed. Obviously, the demand of coordination relative to power and cardiorespiratory endurance is
322 even larger in the star run test than in the 20-m linear sprint test^{31–33} and this could be a reason why
323 the sex effect is smaller for the star run test than for speed.

324 The observed decrease of sex differences in performances with increased demands for coordination
325 seems to be plausible. Coordination draws on many different central (e.g., brain) and peripheral sites

326 (e.g., motor units) within the nervous system³⁴. The more a test engages the brain, the less relevant
327 sex is a performance limiting factor. In summary, the decrease of size of sex effects across tests can
328 most likely be explained by a decrease in the relevance of muscle mass (favouring boys over girls)¹
329 and an increase in the relevance in motor coordination (associated with small or no sex effects)^{15,30}.

330 Within their ninth year of life, older children significantly outperform younger ones in all five fitness
331 tests. This was expected^{3,5,21}, but there were no significant gains for a sample of 240 children in the
332 tenth year of life (see insets of Figure 1). The absence of evidence for nonlinear relations between
333 performance and chronological age as well as the similarity of gains for boys and girls are quite
334 remarkable (see Figure 2), especially in light of what *long-term ontogenetic longitudinal* research
335 reveals when children are a few years older. Obviously, the current test battery does not detect the
336 onset of puberty in the ninth year of life! Nevertheless, the different age-related gains for tests are
337 compatible with what is known about the development of basic physiological parameters: Growth of
338 body mass and height positively influence the performance in *skill-related* physical fitness
339 components of coordination, speed and power. In contrast, increases in body mass have a negative
340 impact on performance in tests such as cardiorespiratory endurance (e.g., 6 min run test) where the
341 own body mass has to be accelerated continuously. The smaller age effect for cardiorespiratory
342 endurance compared to *skill-related* components to the remaining physical fitness components is in
343 line with a study conducted by de Miguel-Etayo et al.³. The authors could not find significant age
344 differences in cardiorespiratory endurance (i.e., 20 m shuttle run test) but amongst other in
345 powerLOW (i.e., standing long jump test) and strength (i.e., absolute handgrip strength test) in a
346 sample of children aged 6 to 10 years. Similar, Viru et al.³⁵ stated that an accelerated improvement in
347 cardiorespiratory endurance occurs at the ages of 11 to 15 years in boys and 11 to 13 years in girls.

348 As far as the differential age effects between the components of physical fitness are concerned, tests
349 of coordination, speed, and powerLOW share the highest correlations among the five tests. These
350 three tests share the relevance of muscle mass yielding power, but they differ in the relevance of
351 coordination. As mentioned above, the sex effect within these three tests (i.e., star run test < 20-m

352 linear sprint test < standing long jump test; see differences between lines in Figure 2) is in line with
353 their ranking on coordination. For age-related gains within the school year, the partial effects yielded
354 star run test > 20-m linear sprint test and standing long jump test, corresponding (roughly) to their
355 ranking on motor coordination. The special status of powerUP (i.e., ball push test) is not only evident
356 with respect to its lower correlations with other tests, but also with respect to the size of the age
357 effect – by far the largest of the five tests (see Figure 2).

358 Obviously, the performance in powerUP was influenced by factors other than physical fitness. We
359 propose that body mass contributes to performance in tests that assess powerUP. Heavier children
360 usually have a higher muscle mass compared to normal weighted children³⁶. While a higher muscle
361 mass positively influences performances in non-weight-bearing tests (e.g., ball push test), it
362 negatively influences performances in weight-bearing tests (e.g., standing long jump test) because
363 the body mass has to be accelerated in contrast to non-weight-bearing tests¹². Therefore, tests for
364 the assessment of powerUP favour heavier children for whom strength of arms and trunk may mask
365 reduced cardiorespiratory endurance and coordination. Thus, a high score for powerUP may be more
366 indicative of a lack of overall physical fitness because it does not measure physical fitness to the
367 same degree than the other four tests (i.e., small correlations with the other four tests and GM). Our
368 results suggest that the best indicator of physical fitness is the average of the first four tests
369 (cardiorespiratory endurance, coordination, speed, and powerLOW). Analyses including body mass
370 (not measured in the present study) could support the interpretation of the special status powerUP
371 in the assessment of physical fitness.

372 The linear mixed model included school as a second random factor, supported the estimation of
373 correlation parameters for test scores and age effects, and revealed a strong correlation ($r = 0.48$)
374 between the age effect and overall physical fitness. The direction of causality is not clear, but in line
375 with Hattie¹⁷. The results are in agreement with the hypothesis that schools differ in how much they
376 promote the development of children's physical fitness. These school differences in the age effect of

377 physical fitness within the third grade were strong enough to yield reliable differences in the overall
378 physical fitness of their children.

379 Our study is not without limitations. First, the five tests do not cover all components of physical
380 fitness¹⁶ and there are alternative tests for each component. Although our results do not necessarily
381 generalize to other components like muscle strength, muscle endurance, or balance¹⁶, they represent
382 those for whom an early detection of puberty was most likely. Second, from a cross-sectional study
383 we cannot know whether the linear gains for tests hold at the individual level or are the result of
384 averaging over individual differences in non-linear growth curves. Obviously, high-density monitoring
385 within a year would be desirable, but such longitudinal data are not without their own problems. For
386 example, how would we separate learning effects due to repeated exposure separate from growth?
387 Motivational factors may also play a role. The longitudinal cardiorespiratory endurance data in Figure
388 1 suggest no further growth or even a decline in performance for 12 year old children⁶. This is
389 obviously not in agreement with what we know about the objective development of
390 cardiorespiratory endurance^{2,4,5}. Third, divergence between cross-sectional and longitudinal profiles
391 is usually cultural change (i.e., cohort effects). Data were accumulated from 2011 to 2019, which is
392 long enough for cohort effects to materialize. The variance of between cohort differences in physical
393 fitness was by far the smallest source examined in this study. There was also no evidence for
394 interactions between cohort, age, and sex in *post-hoc* analyses. Fourth, physical fitness is highly
395 related to biological maturity and more mature youth outperform less mature youth in physical
396 fitness³⁷. This study indicates the strongest test of an early detection of puberty-related divergence
397 with physical-fitness tests that we are aware of. We conclude that puberty-related development has
398 not started or is not strong enough yet in eight-year old children to be picked up with physical fitness
399 or some of its components. A joint analysis with anthropometric measures (i.e., body mass, body
400 height) and status of biological maturity (e.g., peak height velocity, secondary sex characteristics) will
401 allow a stronger test of the hypothesis that the onset of puberty can already be detected in the ninth
402 year of life.

403 To sum up, the *short-term ontogenetic* results of this cross-sectional study revealed test-specific age
404 and sex effects, but no interaction between age and sex despite an abundance of statistical power.
405 According to Ortega et al.²¹ physical fitness data of an individual should be compared with reference
406 values of a sex and age-matched similar general population. Such norms are usually only available on
407 annual or semi-annual basis. Our results suggest that physical education teachers, coaches, or
408 researchers can use a proportional adjustment to adequately evaluate physical fitness of prepubertal
409 school-aged children. Furthermore, especially muscle strength / powerUP should be promoted in
410 sport activities for girls in order to reduce the large sex difference between boys and girls. Lastly,
411 physical education teachers, coaches, or researchers should be careful in the interpretation of the
412 ball push test and should consider that it does not measure physical fitness to the same degree than
413 the other tests.

414 **Methods**

415 ***Sample and study design***

416 This cross-sectional study is part of the ongoing EMOTIKON research project mandated and approved
417 by the Ministry of Education, Youth and Sport of the Federal State of Brandenburg, Germany.

418 Physical fitness tests were carried out by physical education teachers during regular school hours.

419 The Brandenburg School Law requires that parents are comprehensively informed prior to the start
420 of the study. Consent is not needed given that tests are obligatory for both children and schools.

421 Physical fitness tests were administered to *all* third-graders in the state annually at the beginning of
422 the first school semester from 2011 to 2019. Physical fitness tests were also administered to 2009
423 and 2010 cohorts, but at the beginning of the second school semester. Due to the seasonal variation
424 in physical fitness these data were not included.

425 We started with data from 144,045 children. Of those, we included only healthy children who had
426 been enrolled within the legal key date of the Federal State of Brandenburg, that is in a given year of

427 school enrolment they were at least 6.00 and at most 6.99 years old on September 30th and,
428 therefore, varied between 8.00 and 8.99 years in the third grade (n = 110,669). In addition to early-
429 entry (n = 2,664), late-entry (n = 30,457) and children without information about birthdate (n = 255),
430 we did not include children with signs of emotional (e.g., autism) and/or physical disorders (e.g.,
431 disabilities like infantile cerebral palsy) that were evaluated by the responsible and experienced
432 physical education teacher based on a medical clearance (n = 28). After the first iteration, the LMM-
433 based conditional means of the random effects identified one school as an extreme outlier on several
434 tests across the years and the 171 children of that school were excluded as well. Finally, we applied a
435 +/-3 SD criterion to individual test scores which led to the exclusion of another 2,175 children (2%).
436 This left us with 108,295 children (i.e., 75% of those tested) from 515 different schools.

437 ***Physical fitness tests***

438 Physical fitness was assessed with the EMOTIKON test battery. The five tests measured
439 cardiorespiratory endurance (i.e., 6 min run test), coordination (i.e., star run test), speed (i.e., 20-m
440 linear sprint test), power of lower limbs (powerLOW [i.e., standing long jump test]), and power of
441 upper limbs (powerUP [i.e., ball push test]). The EMOTIKON test battery officially includes six tests.
442 Up to 2012 the sixth test was the stand and reach test (flexibility) that was then exchanged against
443 the single-leg balance test (balance). Due to the much smaller number of scores and their
444 confound with cohort these tests were not included in the analyses. The five tests yielded 525,126
445 scores from the 108,295 children (i.e., 3% missing test scores).

446 Qualified physical education teachers of each school administered the tests according to
447 standardized test protocols during the regular physical education classes in the participating schools
448 (www.uni-potsdam.de/en/emotikon/projekt/methodik for further information on the test
449 protocols). Teachers were instructed in a standardized assessment through an advanced training.
450 Tests were always conducted in the morning between 8 and 12 o'clock. Encouragement to achieve
451 the best performance was permitted. Teachers carried out individualized warm-ups before testing.

452 *Cardiorespiratory endurance*

453 Cardiorespiratory endurance was assessed with the 6 min run test. Children had to run as far as they
454 could within six minutes around an official volleyball field (9 x 18 m, every 9 m a pylon/marker was
455 set beside the running court [i.e., six pylons around the field]) at a self-paced velocity. Split time was
456 given every minute. The maximal distance achieved during the six minutes in meters to the nearest
457 nine-meters marker was used as dependent variable in the analysis. The 6 min run test was reliable
458 (test-retest) in children aged 7 to 11 years with an intraclass correlation coefficient (ICC) of 0.92³⁸.

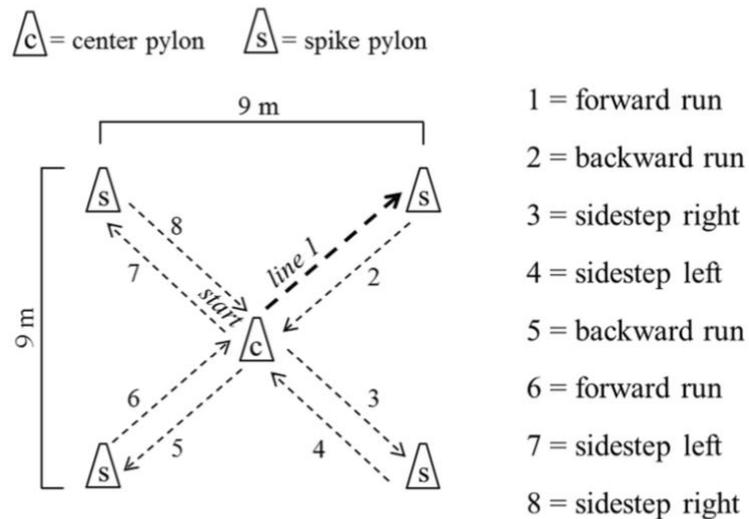
459 *Coordination*

460 Coordination under time pressure was tested with the star run test (see Figure 5). Children had to
461 complete a parkour with different movement directions and movement forms (i.e., running forward,
462 running backward, side-steps to the left side, side-steps to the right side). The parkour had to be
463 performed in a given order over a 9 x 9 m star-shaped area where each of the four spikes is marked
464 by a pylon. After starting in the centre of the star, children had to complete the parkour as fast as
465 they could by running in every movement form two times within the given sequence. They had to
466 touch each pylon with the hand. The whole covered distance is 50.912 m. The faster of two test trials
467 was used in the analysis. The shortest time for completing the parkour in seconds to the nearest 1/10
468 second was measured using a stopwatch and was used as dependent variable in the analysis. The
469 star run test was reliable (test-retest) in 8 to 10 year old children with an ICC of 0.68³⁹.

470 *Speed*

471 Speed was assessed with the 20-m linear sprint test. After an acoustic signal, the children
472 had to sprint out of a frontal erect posture as fast as they could over a distance of 20 m for
473 two times; the faster test trial was used in the analysis. The shortest time for sprinting the 20
474 m in seconds to the nearest 1/10 second was measured using a stopwatch and was used as

475 dependent variable in the analysis. The 20-m linear sprint test was reliable (test-retest) in
 476 children aged 7 to 11 years with an ICC of 0.90³⁸.



477
 478 **Figure 5.** Schematic description of the star run test (adapted from Golle et al.⁶).

479 *Power of lower limbs (PowerLOW)*

480 PowerLOW was tested using the standing long jump test. Out of a standing frontal posture
 481 the children had to jump as far as they could. The participants had to land with both feet
 482 together. They were allowed to swing their arms prior to and during the jump, but after
 483 landing the hands were not allowed to touch the floor. The distance in meters to the nearest
 484 one centimeter between toes at take-off and heels at landing was determined using a
 485 measuring tape; the better of two test trials was used in the analysis. The standing long jump
 486 test was reliable (test-retest) in children aged 6 to 12 years with an ICC of 0.94⁴⁰.

487 *Power of upper limbs (PowerUP)*

488 PowerUP was assessed through the ball push test. From a standing position the children had
 489 to push a 1 kg medicine ball starting in front of the chest with both hands as far as they

490 could for two times; the better of two test trials of longest pushing distance was used in the
491 analysis. The maximal ball push distance in meters to the nearest ten centimeters was
492 determined with a measuring tape and used as dependent variable in the analysis. The ball
493 push test was reliable (test-retest) in children aged 8 to 10 years with an ICC of 0.81³⁹.

494 **Statistics**

495 Pre- and post-processing of data were carried out in the R environment of statistical computing⁴¹
496 using the *tidyverse* package⁴². For measures of cardiorespiratory endurance (i.e., 6 min run test),
497 powerLOW (i.e., standing long jump test) and powerUP (i.e., ball push test), higher scores indicated
498 better physical fitness. For measures of coordination (i.e., star run test) and speed (i.e., 20-m linear
499 sprint test), a Box-Cox distributional analyses indicated that a reciprocal transformation brought
500 scores in line with the assumption of a normal distribution⁴³. Therefore, we converted scores from
501 seconds to meters/seconds (i.e., pace scores; star run test = 50.912 [m] / time [s]; 20-m linear sprint
502 test = 20 [m] / time [s]). These transformations also had the advantage that a large value was
503 indicative of a good physical fitness for all five measures.

504 For each test, we determined the ± 3 SD boundary separately for boys and girls. Measurement
505 outside these boundaries were usually implausible (i.e., recording errors) or extreme outliers. They
506 were treated as missing values (3%). Finally, we converted scores within tests (aggregated over boys
507 and girls) to z-scores to facilitate comparison of test, age and sex effects.

508 Statistical inference was based on a linear mixed model (LMM) estimated with the *MixedModels*
509 package⁴⁴ in the *Julia* programming language⁴⁵. The LMM included child (N = 108,295), school (N =
510 515), and cohort (N = 9) as three random factors; the total number of observations (i.e., max = 5 per
511 child) was 525,126.

512 As fixed effects, we specified four sequential-difference contrasts for the five tests: (H1) coordination
513 vs. cardiorespiratory endurance, (H2) speed vs. coordination, (H3) powerLOW vs. speed, and (H4)
514 powerUP vs. powerLOW. Also included were the effect of age (centered at 8.5 years) as a second-

515 order polynomial trend, the effect of sex (boys – girls), and all interactions between contrasts, age,
516 and sex. Given the large number of observations, children, and schools, we adopted a two-sided z-
517 value > 3.0 as significance criterion for the interpretation of fixed effects.

518 Child, school, and cohort were included as random factors. With three random factors there was a
519 need for selecting a random-effect structure that included theoretically relevant and reliable
520 variance components (VCs) and correlation parameters (CPs), but was also still supported by the data
521 (i.e., was not overparameterized). Tests varied within children, schools, and cohorts; age and sex
522 varied between children, but within schools and within cohorts. Therefore, in principle, VCs and CPs
523 of linear effects of age and sex could be estimated for schools and cohorts, but not for children.

524 Parsimonious model selection occurred in two major steps without knowledge or consideration of
525 fixed-effect estimates⁴⁶; details are provided in Supplement A. We started with a model including
526 Grand Mean (varying intercepts) for all three random factors and, given the large numbers of
527 108,926 children and 515 schools and the small number of nine cohorts, included also test-related
528 VCs and CPs for child and school and age-related and sex-related VCs and CPs for school, but not for
529 cohort. This LMM *m1* was well supported by the data. In the second major step, we increased the
530 complexity of the random-effect structure for cohort by adding test-related VCs (LMM *m2*), then
531 test-related CPs (LMM *m3*), and finally age- and sex-related VCs and CPs (LMM *m4*).

532 LMM *m4* was not supported by the data (i.e., the fit was singular) and did not significantly improve
533 the goodness of fit over LMM *m3*; $\Delta \chi^2 (13) = 14.11, p = 0.37$. LMM *m3* improved the goodness of
534 fit over LMM *m2* according to the likelihood ratio test, $\chi^2 (10) = 48.45, p < 0.001$, but not when the
535 increase in model complexity is penalized according to BIC (i.e., LMM *m2* = 1.27609e6 and LMM *m3* =
536 1.27617e6). As we had no directed hypotheses relating to test-related CPs for the factor cohort, we
537 stayed with LMM *m2* which represented a very large improvement in goodness of fit relative to LMM
538 *m1*; $\chi^2 (4) = 1489.57, p < 0.001$. We also estimated LMM *m2* with two alternative parameterizations
539 that did not change the goodness of fit, but yielded information about CPs between test scores
540 instead of test effects (i.e., contrasts). Finally, we fitted two control LMMs to test the significance of

541 quadratic age trends for fixed effects and the absence of evidence for sex x age interactions
542 separately for each fitness component (i.e., nested within the five levels of the factor test).

543 **References**

- 544
- 545 1. Malina, R. M., Bouchard, C. & Bar-Or, O. *Growth, Maturation, and Physical Activity*. (Human
- 546 Kinetics, 2004).
- 547 2. Tomkinson, G. R. *et al.* European normative values for physical fitness in children and
- 548 adolescents aged 9-17 years: Results from 2 779 165 Eurofit performances representing 30
- 549 countries. *Br. J. Sports Med.* **52**, 1445–1456 (2018).
- 550 3. De Miguel-Etayo, P. *et al.* Physical fitness reference standards in European children: The
- 551 IDEFICS study. *Int. J. Obes.* **38**, S57–S66 (2014).
- 552 4. Santos, R. *et al.* Physical fitness percentiles for Portuguese children and adolescents aged 10-
- 553 18 years. *J. Sports Sci.* **32**, 1510–1518 (2014).
- 554 5. Tambalis, K. D. *et al.* Physical fitness normative values for 6–18-year-old Greek boys and girls,
- 555 using the empirical distribution and the lambda, mu, and sigma statistical method. *Eur. J.*
- 556 *Sport Sci.* **16**, 736–746 (2016).
- 557 6. Golle, K., Muehlbauer, T., Wick, D. & Granacher, U. Physical fitness percentiles of german
- 558 children aged 9-12 Years: Findings from a longitudinal study. *PLoS One* **10**, 9–11 (2015).
- 559 7. Niessner, C. *et al.* Representative Percentile Curves of Physical Fitness From Early Childhood
- 560 to Early Adulthood: The MoMo Study. *Front. Public Heal.* **8**, 1–9 (2020).
- 561 8. Stratton, G. & Oliver, J. L. The impact of growth and maturation on physical performance. in
- 562 *Strength and conditioning for young athletes: science and application* (eds. Lloyd, R. S. &
- 563 Oliver, J. L.) 3–20 (Routledge, 2020).
- 564 9. Marshall, W. A. & Tanner, J. M. Variations in pattern of pubertal changes in girls. *Arch. Dis.*
- 565 *Child.* **44**, 291–303 (1969).
- 566 10. Rosenfield, R. L., B.Lipton, R. & Drum, M. L. Thelarche, pubarche, and menarche attainment in
- 567 children with normal and elevated body mass index. *Pediatrics* **123**, 84–88 (2009).
- 568 11. Marshall, W. A. & Tanner, J. M. Variations in pattern of pubertal changes in boys. *Arch. Dis.*
- 569 *Child.* **45**, 13–23 (1970).

- 570 12. Drenowatz, C., Hinterkörner, F. & Greier, K. Physical Fitness and motor competence in upper
571 Austrian elementary school children — study protocol and preliminary findings of a state-wide
572 fitness testing program. *Front. Sport. Act. Living* **3**, 1–11 (2021).
- 573 13. Kanehisa, H., Yata, H., Ikegawa, S. & Fukunaga, T. A cross-sectional study of the size and
574 strength of the lower leg muscles during growth. *Eur. J. Appl. Physiol. Occup. Physiol.* **72**, 150–
575 156 (1995).
- 576 14. Altman, D. G. & Bland, J. M. Absence of evidence is not evidence of absence. *Bmj* **311**, 485
577 (1995).
- 578 15. Overman, W. H. Sex differences in early childhood, adolescence, and adulthood on cognitive
579 tasks that rely on orbital prefrontal cortex. *Brain Cogn.* **55**, 134–147 (2004).
- 580 16. Caspersen, C., Powell, K. & Christenson, G. Physical Activity, Exercise, and Physical Fitness:
581 Definitions and Distinctions for Health-Related Research CARL. *Notes Queries* **100, No. 2**, 125–
582 131 (1985).
- 583 17. Hattie, J. A. C. *Visible learning: A synthesis of over 800 meta-analyses relating to achievement*.
584 (Routledge, 2009).
- 585 18. Fühner, T., Kliegl, R., Arntz, F., Kriemler, S. & Granacher, U. An Update on Secular Trends in
586 Physical Fitness of Children and Adolescents from 1972 to 2015: A Systematic Review. *Sport.*
587 *Med.* **51**, 303–320 (2021).
- 588 19. Tomkinson, G. R., Lang, J. J. & Tremblay, M. S. Temporal trends in the cardiorespiratory fitness
589 of children and adolescents representing 19 high-income and upper middle-income countries
590 between 1981 and 2014. *Br. J. Sports Med.* **53**, 478–486 (2019).
- 591 20. Spengler, S., Rabel, M., Kuritz, A. M. & Mess, F. Trends in motor performance of first graders:
592 A comparison of cohorts from 2006 to 2015. *Front. Pediatr.* **5**, 1–7 (2017).
- 593 21. Ortega, F. B. *et al.* Physical fitness levels among European adolescents: The HELENA study. *Br.*
594 *J. Sports Med.* **45**, 20–29 (2011).
- 595 22. Woll, A., Kurth, B. M., Opper, E., Worth, A. & Bös, K. The ‘Motorik-Modul’ (MoMo): Physical

- 596 fitness and physical activity in German children and adolescents. *Eur. J. Pediatr.* **170**, 1129–
597 1142 (2011).
- 598 23. Catley, M. J. & Tomkinson, G. R. Normative health-related fitness values for children: Analysis
599 of 85347 test results on 9-17-year-old Australians since 1985. *Br. J. Sports Med.* **47**, 98–108
600 (2013).
- 601 24. Roriz de Oliveira, M. S., Seabra, A., Freitas, D., Eisenmann, J. C. & Maia, J. Physical fitness
602 percentile charts for children aged 6-10 from Portugal. *J. Sports Med. Phys. Fitness* **54**, 780–
603 792 (2014).
- 604 25. Ramos-Sepúlveda, J. A., Ramírez-Vélez, R., Correa-Bautista, J. E., Izquierdo, M. & García-
605 Hermoso, A. Physical fitness and anthropometric normative values among Colombian-Indian
606 schoolchildren. *BMC Public Health* **16**, 1–15 (2016).
- 607 26. Beunen, G. & Thomis, M. Muscular strength development in children and adolescents.
608 *Pediatr. Exerc. Sci.* **12**, 174–197 (2000).
- 609 27. Round, J. M., Jones, D. A., Honour, J. W. & Nevill, A. M. Hormonal factors in the development
610 of differences in strength between boys and girls during adolescence: A longitudinal study.
611 *Ann. Hum. Biol.* **26**, 49–62 (1999).
- 612 28. Armstrong, N. & Welsman, J. Aerobic fitness: What are we measuring? *Med. Sport Sci.* **50**, 5–
613 25 (2007).
- 614 29. Haywood, K. M. & Getchell, N. *Life span motor development*. (Human Kinetics, 2009).
615 doi:10.1016/0167-9457(87)90023-6.
- 616 30. Ardila, A., Rosselli, M., Matute, E. & Inozemtseva, O. Gender Differences in Cognitive
617 Development. *Dev. Psychol.* **47**, 984–990 (2011).
- 618 31. Ludyga, S., Gerber, M., Pühse, U., Looser, V. N. & Kamijo, K. Systematic review and meta-
619 analysis investigating moderators of long-term effects of exercise on cognition in healthy
620 individuals. *Nat. Hum. Behav.* **4**, 603–612 (2020).
- 621 32. Schmidt, M. *et al.* Disentangling the relationship between children’s motor ability, executive

- 622 function and academic achievement. *PLoS One* **12**, (2017).
- 623 33. Koutsandréou, F., Wegner, M., Niemann, C. & Budde, H. Effects of motor versus
624 cardiovascular exercise training on children's working memory. *Med. Sci. Sports Exerc.* **48**,
625 1144–1152 (2016).
- 626 34. Viru, A. *et al.* Critical Periods in the Development of Performance Capacity During Childhood
627 and Adolescence. *Eur. J. Phys. Educ.* **4**, 75–119 (1999).
- 628 35. Viru, A. *et al.* Age periods of accelerated improvement of muscle strength, power, speed and
629 endurance in the age interval 6-18 years. *Biol. Sport* **15**, 211–227 (1998).
- 630 36. Ducher, G. *et al.* Overweight children have a greater proportion of fat mass relative to muscle
631 mass in the upper limbs than in the lower limbs: Implications for bone strength at the distal
632 forearm. *Am. J. Clin. Nutr.* **90**, 1104–1111 (2009).
- 633 37. Jones, M. A., Hitchen, P. J. & Stratton, G. The importance of considering biological maturity
634 when assessing physical fitness measures in girls and boys aged 10 to 16 years. *Ann. Hum.*
635 *Biol.* **27**, 57–65 (2000).
- 636 38. Bös, K. *Deutscher Motorik-Test 6-18*. (Czwalina, 2009).
- 637 39. Schulz, S. *The reliability of the Star coordination run and the 1-kg medicine ball push - physical*
638 *fitness tests used in the EMOTIKON-study*. (University of Potsdam, 2013).
- 639 40. Fernandez-Santos, J. R., Ruiz, J. R., Cohen, D. D., Gonzalez-Montesinos, J. L. & Castro-Pinero, J.
640 Reliability and validity of tests to assess lower-body muscular power in children. *J. Strength*
641 *Cond. Res.* **29**, 2277–2285 (2015).
- 642 41. Team, R. C. R a language and environment for statistical computing. [https://www.r-](https://www.r-project.org/)
643 [project.org/](https://www.r-project.org/) (2020).
- 644 42. Wickham, H. *et al.* Welcome to the Tidyverse. *J. Open Source Softw.* **4**, 1686 (2019).
- 645 43. Box, G. E. P. & Cox, D. R. An analysis of transformations. *J. R. Stat. Soc. Ser. B* **26**, 211–243
646 (1964).
- 647 44. Bates, D. *et al.* JuliaStats/MixedModels.jl:v3.4.0. vol. 1 20071787–20071787

- 648 <https://doi.org/10.5281/zenodo.4589101> (2021).
- 649 45. Bezanson, J., Edelman, A., Karpinski, S. & Shah, V. B. Julia: A fresh approach to numerical
650 computing. *SIAM Rev.* **59**, 65–98 (2017).
- 651 46. Bates, D., Kliegl, R., Vasishth, S. & Baayen, H. Parsimonious Mixed Models. *ArXiv:1506.04967*
652 <http://arxiv.org/abs/1506.04967> (2015).
- 653

654 **Acknowledgements**

655 The study was commissioned and funded by the Ministry of Education, Youth, and Sport of the
656 Federal State of Brandenburg, Germany. The funders had no role in study design, data collection and
657 analysis, decision to publish, or preparation of the manuscript. Reinhold Kliegl was supported by the
658 Center for Interdisciplinary Research, Bielefeld (ZiF)/Cooperation Group "Statistical models for
659 psychological and linguistic data".

660 **Author contributions**

661 TF, KG, and UG contributed to conception and design; TF and KG organized data collection; RK and TF
662 carried out data analysis; TF and RK wrote the first draft of the manuscript and all authors were
663 involved in iterative revisions; all authors provided final approval of the version to be published and
664 agreed to be accountable for all aspects of the work.

665

666 **Additional Information**

667 **Competing Interests Statement**

668 Thea Fühner, Urs Granacher, Kathleen Golle, and Reinhold Kliegl declare that they have no conflicts
669 of interest relevant to the content of this cross-sectional study.

670 **Availability of material, data, and code**

671 Data, Julia, and R scripts are available at an Open Science Framework (OSF) repository:
672 <https://osf.io/2d8rj/>

Figures

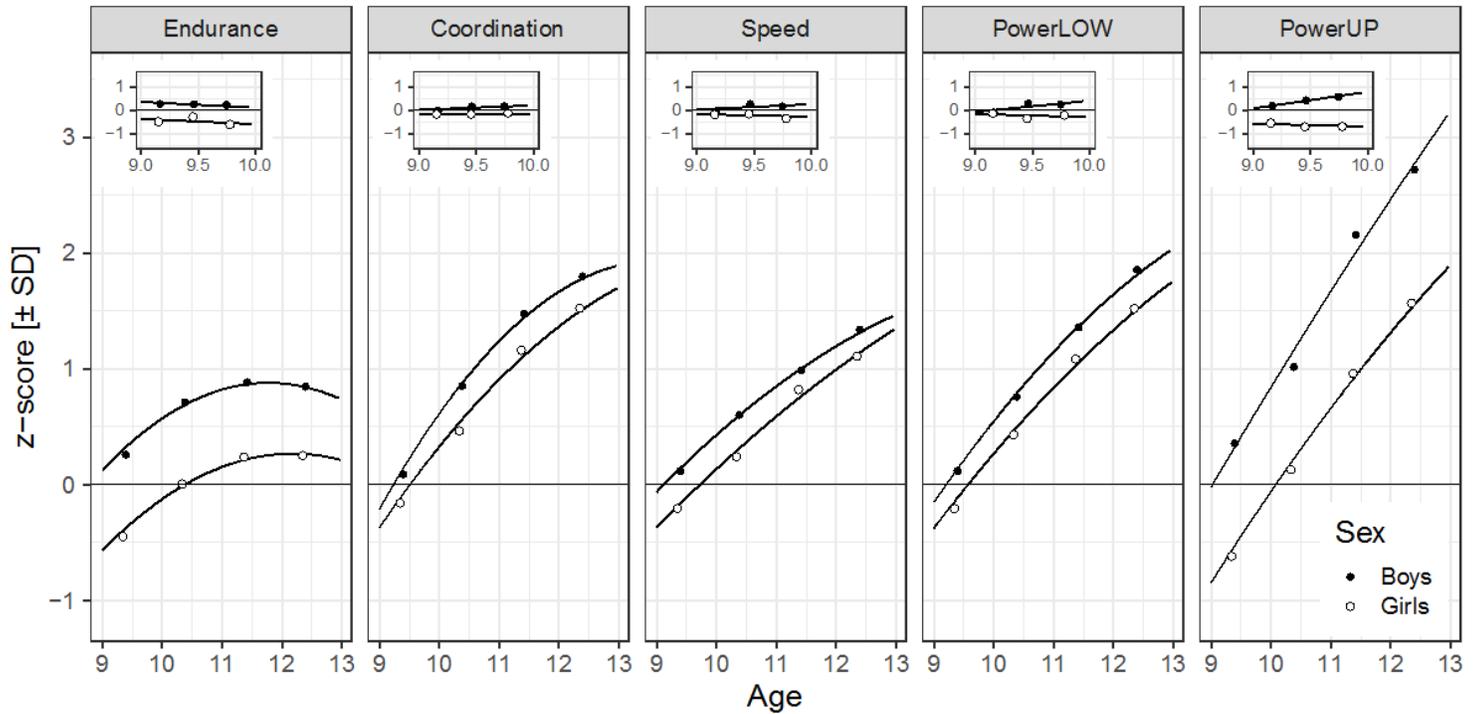


Figure 1

Physical fitness curves of a longitudinal sample of 240 German boys (closed circles) and girls (open circles) followed from age 9 to 12 years for endurance = cardiorespiratory endurance (i.e., 9 min run test), coordination (i.e., running in a star like pattern), speed (i.e., 50-m linear sprint test), powerLOW = power of lower limbs (i.e., triple hop test), and powerUP = power of upper limbs (i.e., ball push test). The insets show for each test score the regression on age for the first assessment when children were between 9.00 and 9.99 years old. Also shown are the means for groups of boys and girls binned into three age groups (i.e., 9.00-9.33; 9.34-9.66; 9.67-9.99). Error bands are 95% Cis. Data are from Golle et al.6.

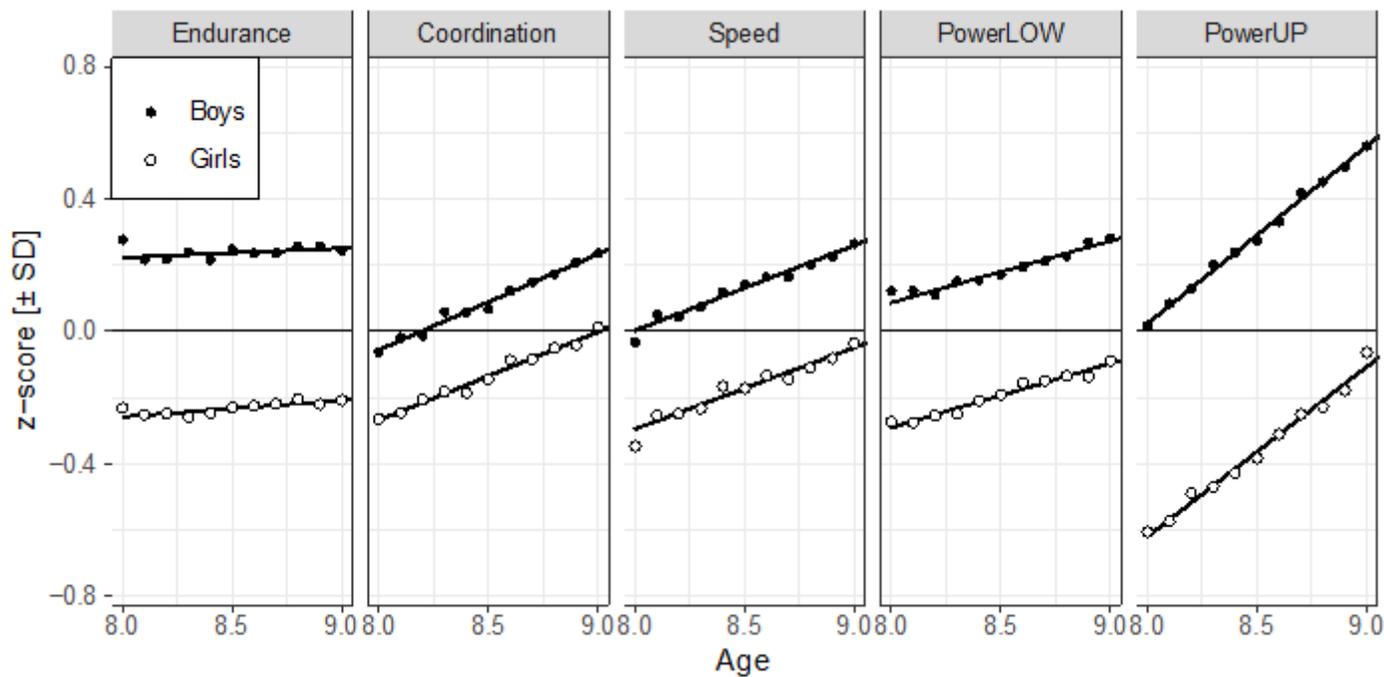


Figure 2

Performance differences between 8.0 and 9.0 years by sex in the five physical fitness tests presented as z-transformed data computed separately for each test. Endurance = cardiorespiratory endurance (i.e., 6 min run test), Coordination = star run test, Speed = 20-m linear sprint test, PowerLOW = power of lower limbs (i.e., standing long jump test), PowerUP = power of upper limbs (i.e., ball push test), SD = standard deviation. Points are binned observed child means; lines are simple regression fits to the observations; 95% confidence intervals for means ≈ 0.05 are not visible.

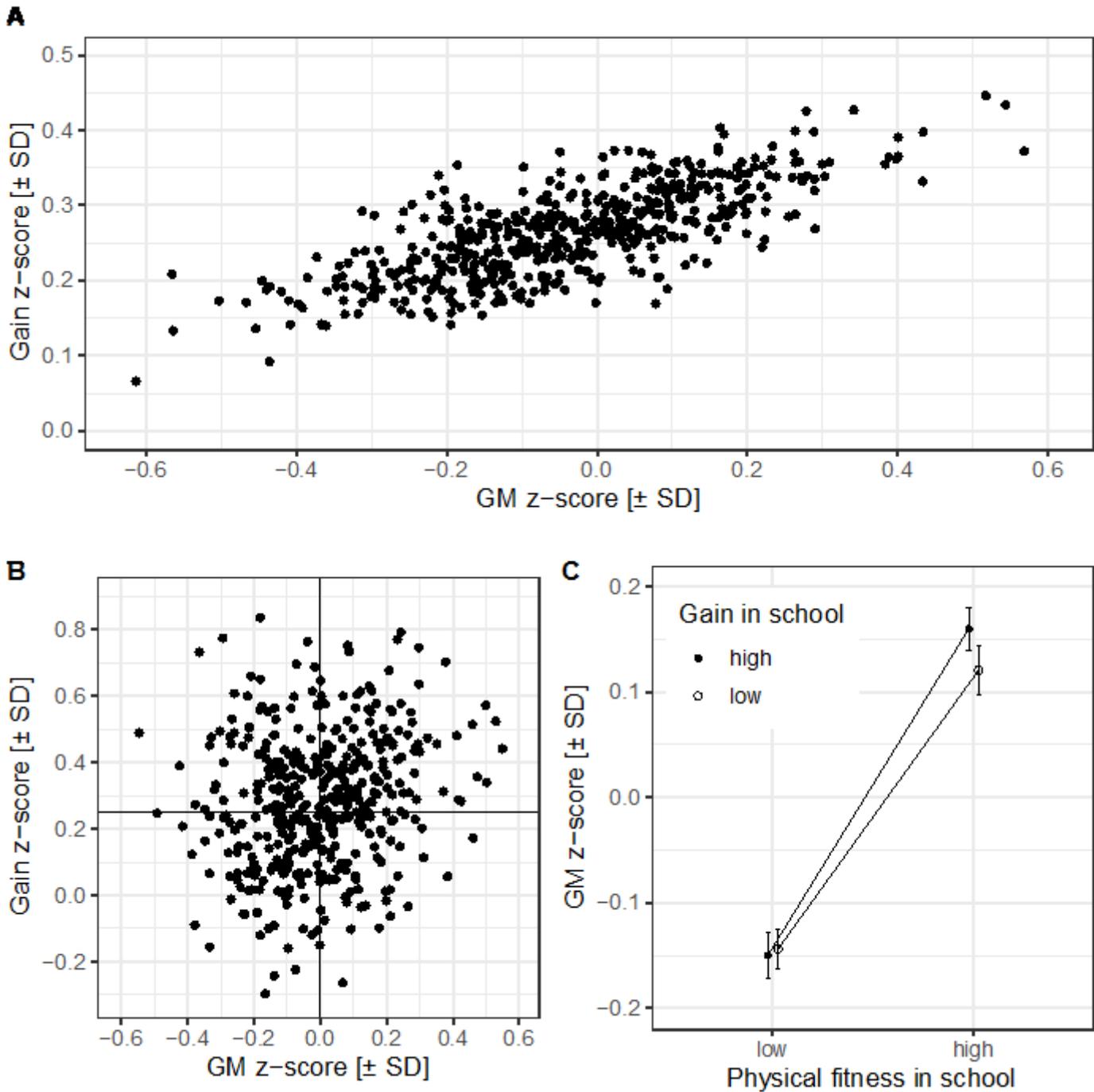


Figure 3

(a) Visualization of correlation parameter between Grand Mean (GM) and yearly gains (age effect) of 515 schools using conditional means resulting from shrinkage correction of observed data with LMM parameters. (b) Scatterplot of observed GMs and yearly gains (within-school age slopes) for 421 schools reporting data from more than 30 boys and 30 girls. (c) After sorting schools into groups of high vs. low physical fitness (split at 0) and high vs. small gain-rate (split at the gain of +0.25), the interaction corresponding to the CP becomes visible; error bars are 95% CIs. SD = standard deviation.

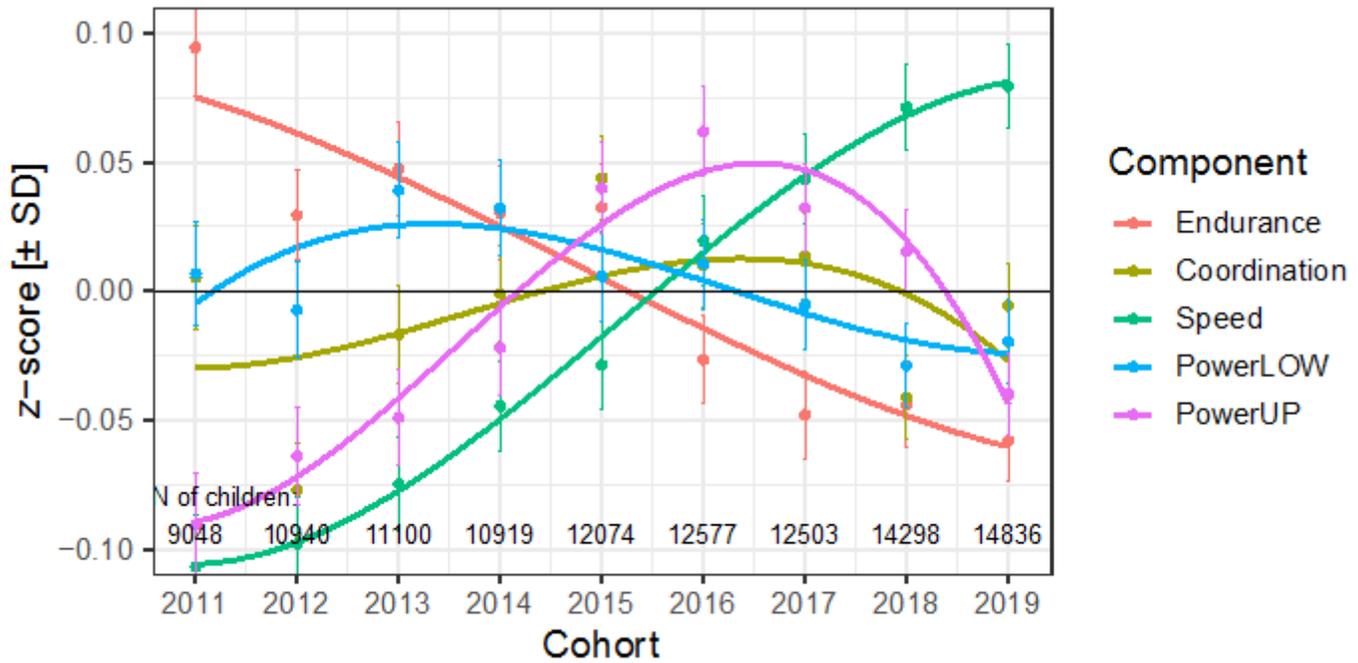


Figure 4

Cohort-related change of components of physical fitness. Points are observed means with 95% CIs. Lines are third-order polynomial trends fitted to children's scores along with 95% error. Note the much smaller range of the y-axis (i.e., from -0.10 to +0.10) compared to age effects shown in Figure 2. Endurance = cardiorespiratory endurance (i.e., 6 min run test), Coordination = star run test, Speed = 20-m linear sprint test, PowerLOW = power of lower limbs (i.e., standing long jump test), PowerUP = power of upper limbs (i.e., ball push test).

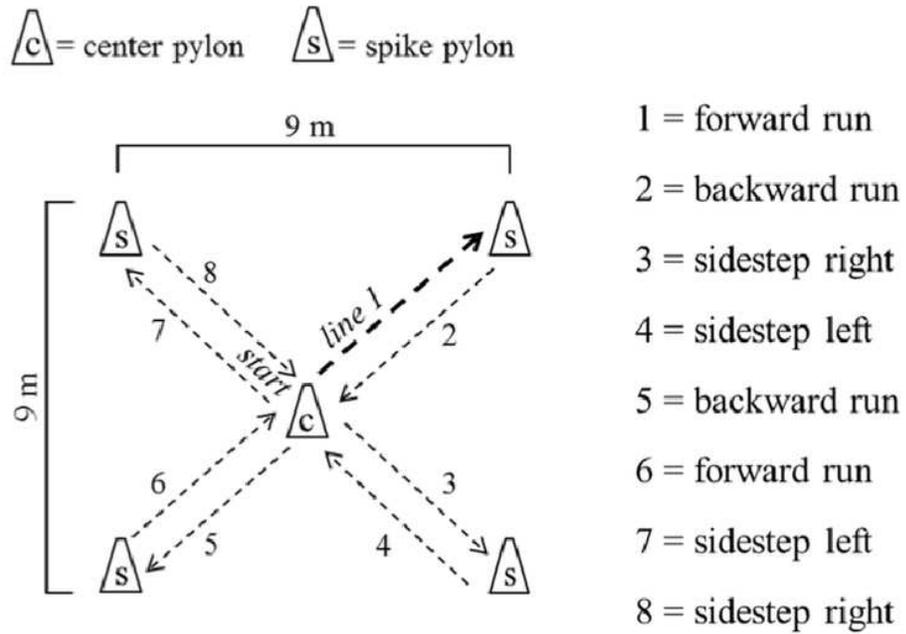


Figure 5

Schematic description of the star run test (adapted from Golle et al.6).

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [SupplementFGGK21.pdf](#)