

13 **Abstract:**

14 Despite featuring extreme physiological adaptations integration of wildlife species into the
15 modern ‘omics’ frameworks are limited due to the sparsity in the data. To address the sparsity
16 limitation a self-tuning sliding-window framework was developed for the identification of the
17 regional poly-CpG methylation architecture associated with phenotypic traits. Under the
18 framework the iteratively expanded window aligns against each chromosome for two
19 constraints: (i) the window must increase its penalised R^2 statistic derived from univariate
20 (ADVI-EWAS) and multivariate (Horseshoe) CpG-phenotype models and (ii) both models but
21 agree upon the direction of the methylome-phenotype association. Any windows not passing
22 the two conditions were discarded if they lacked sufficient CpG density, otherwise were
23 retained but ceased to grow. Application of the framework to the blood methylome of 69
24 California Sea Lions (*Zalophus californianus*) detected three ‘trusted’ windows across the
25 entire genome. Each of the windows were uniquely mapped to a single gene, these genes were
26 GCSH, DNMT1, LMX1A. All three genes are mechanistically linked to the phenotypic trait
27 and are consistent with the extreme phenotype of marine mammals. The results were
28 independent of the age of the individuals tested. The results were not predictive of phenotype
29 when applying epigenome wide association studies (EWAS); however, application of the poly-
30 CpG architecture provided a phenotypically linked signal that was biologically interpretable.
31 The sliding window method is computationally efficient, isolates short, directionally stable
32 genome regions and outperforms single-site analyses. The generalisable statistical framework
33 provided within the study could extract phenotype-linked molecular architecture from sparse
34 wildlife methylomes.

35 Key words – Physiology, Phenotype, Methylation, Sliding-window

36 **1. Introduction:**

37 All 'omics' data is inherently sparse, often requiring large sample sizes for inferencing
38 (1, 2). The application of 'omics' to wildlife species is even more pertinent, with smaller sample
39 sizes afforded due to the opportunistic nature of sampling and greater variability present (3, 4).
40 Wildlife species are frequently neglected within the 'omics' (3). Methylation studies are just
41 one example of the phenomenon – with the inherent environmental stressors further
42 complicating the consideration of the phenotype (5).

43 Unlike with genomics, methylation studies face greater variability within the models as
44 a result of the phenotypic influences upon the methylome and the methylomes influences upon
45 the phenotype (5). As such, wildlife methylome studies are typically restricted to molecular
46 aging rather than phenotype-methylation association studies (6, 7). Single CpG models –
47 including epigenome wide association studies (EWAS) – are just as sparse, if not more, sparse
48 than genome wide association studies (GWAS). Both frameworks are already sparse in wildlife
49 and frequently return an absence of significance in well studies model organisms (1, 2, 8).
50 However, the absence of a causal variable is not sufficient to declare the absence of evidence
51 within biological models (1, 2). Improving the sample size will reduce sparsity and enable
52 greater significance, this is not always achievable for wildlife due to limited sampling
53 opportunities (1).

54 Regional poly-CpG architecture dictates that the cumulative effect of CpGs within a
55 defined genomic region is more biologically meaningful than that of a single-CpG model of
56 phenotypic associations (9). Under this hypothesis, the accumulation of methylated sites
57 results in the refolding of the chromosome – further reducing the expression of the particular
58 chromosomal regions (9). The framework must tackle distance and issues of capturing noise in
59 the poly-CpG model.

60 Under the sliding window framework the genome is regionally scanned across individual
61 chromosomes and for each CpG iteratively expanded within the chromosome for the
62 optimisation of a penalised R^2 parameter. The expansion of the window will continue until the
63 either the R^2 penalised value doesn't increase sufficiently for each inclusion of a CpG and the
64 agreement of a biologically meaningful value for the window's collective CpG scores and their
65 estimated effects on the phenotype. Any windows failing to meet these conditions will be
66 terminated and those under the threshold of CpG inclusion number will be permanently
67 discarded. Application of a sliding window model reduces the computational burden and adds
68 biological meaning (regionality) to the framework in a similar manner to GWAS (8).

69 Demonstration of the sliding window framework is necessary. The extreme physiology
70 of the California Sea Lion (*Zalophus californianus*) renders the species suitable to study.
71 Recently *Z. californianus* was molecularly aged through the methylome (6). Although these
72 results found a strong association between the sex and chronological age upon the epigenetic
73 age of the individuals was determined there was a lack of phenotypic inclusion in the model
74 beyond these two variables. Thus, these samples demonstrate untapped potential. The
75 maximum aerobic dive capacity was chosen as an ideal phenotype for the identification of the
76 genes associative with a taxa-specific phenotype as the variable was a linear function of age –
77 Sea Lions experience increasing physiological demands as they mature into independently
78 feeding individuals – a process that can take years to complete (10, 11). It was hypothesised
79 that the increasing physiological demands mid-dive is a driver for methylation changes – which
80 will be measured through the sliding window framework (6).

81 The aim of the study is to apply a self-tuning sliding window framework to the analysis
82 of the phenotype-linked poly-CpG structuring of the methylome of 69 individual blood samples
83 of *Z. californianus*.

84 **Methods:**

85 **2.1 Methods overview:**

86 **2.1.1 Introduction to the methods**

87 An eight-stack pipeline was used to analyse the relationship between the methylome and
88 the phenotype solely using Python-3.13. Statements of the full mathematical formulae to
89 reproduce the statistics are listed in the Supplementary Methods. All code is available in the
90 GitHub repository for the study ([tsto3616/poly-CpG: The repository encodes the sliding
91 window pipeline, accompanying statistics and data for a study on the California Sea Lion's
92 epigenome](#)).

93 All data included within the study was collected from open-access sources and the
94 sources of the data both cite ethics approval (6, 10). No new animal procedures were conducted
95 for these experiments.

96 **2.1.2 Methylation data**

97 Robeck, Haghani (6) studied the methylation of multiple pinniped species, of these
98 samples 69 blood samples of *Z. californianus* were used for this study. The samples represented
99 multiple age groups and the normalised methylation data was sourced from Gene Expression
100 Omnibus (GEO accession: GSE227319). Methylation annotations were collected from the
101 GitHub repository featured as part of the study by Robeck, Haghani (6). The genome/ gene
102 annotation file for *Z. californianus* was collected from NCBI (GCA_009762295.2). Handling
103 of the raw methylation data files is not discussed as it was not part of the study. Rather the
104 normalised methylation data from GEO (GSE227319) were processed.

105 **2.1.3 Phenotypic data**

106 The oxygen carrying capabilities of *Z. californianus* were determined through the
107 formula published by Weise and Costa (10). The formula is $MaxADC =$
108 $0.2357(years\ old) + 1.3571$, where the *MaxADC* represents the maximum aerobic diving
109 capacity (in minutes) for the age of the animal. The maximum aerobic diving capacity is treated
110 as a continuous phenotypic variable throughout downstream analyses.

111 **2.1.4 Construction of the master file and basic statistics**

112 The workflow involved the conversion of published β -values to M-scores using the
113 standard logit transformation (Supplementary Methods 1.2) after removing the extreme values
114 (close to 0 or 1) to prevent the production of infinite M-scores (6). The top 10,000 CpG sites
115 for which there was a high rate of variance across the 69 samples were the only CpGs included
116 in this study. Inclusion of greater numbers of CpG sites would create a significant
117 computational burden with minimal benefit. Minimal variation would only weaken the signal
118 to noise ratio and so these CpGs were removed.

119 **2.1.5 Overview of the multi-stack modelling framework**

120 Stack One clustered the methylation profiles of *Z. californianus* according to the M-
121 scores of the methylation data by applying HDBSCAN.

122 Stack Two calculated the latent age for downstream analyses using the epigenomic
123 calculated age and the chronological age. The latent age was used to determine the reliability
124 of the clustering from Stack One.

125 Stack Three predicted the maximum aerobic dive capacity through a global Hierarchical
126 Bayesian model – as determined by the methylation clusters and latent age reliability of the
127 clusters (Stack One and Two respectively).

128 Stack Four predicted the relationship of CpGs through univariate epigenome wide
129 association study (EWAS) sampled with Automatic Differentiation Variation Inference
130 (ADVI). The ADVI model was informed by latent age (Stack Two) and the residuals from
131 Stack Three as the outcome variable. The conversion of the results into the Bayesian
132 probabilities and their standard deviation into measurable outcomes (Bayesian probability).

133 Stack Five involved the determination of a Horseshoe Shrinkage Bayesian model was
134 fitted to effectively handle the sparse CpG data through a multivariate model.

135 Stack Six involved the isolation of CpGs predicted to be a significant predictor of
136 maximum aerobic diving capacity for both the ADVI and Horseshoe Bayesian shrinkage
137 model. This involved the mapping of the CpGs to the genomic coordinates.

138 Stack Seven evaluated the role of regional context in the poly-CpG scores for which a
139 sliding window regression task was fitted. This model involved a moveable window for which
140 the poly-CpG scores were considered for the maximisation of the correlation while limiting the
141 sliding windows to the one chromosome to enable biologically relevant signal generation. The
142 sliding windows were overlaid onto the genes of *Z. californianus* for the isolation of genes
143 for which there is an association with the maximum aerobic dive capacity.

144 Stack Eight produced a comparison of the earlier models using trajectory analysis for the
145 determination of the methylation profiling in response to the models and their variables. The
146 pipeline involved dimension reduction through PCA, k-NN clustering, MST and DPT analyses
147 for the determination of the pseudotime for which the methylation profiles differ.

148 **2.2 Stack One — High-variance CpG selection, PCA reduction, and HDBSCAN** 149 **clustering**

150 Stack One performed clustering of the genome-wide methylation profiles for the
151 identification of the methylation architectural structure across individual California Sea Lions.

152 The first 20 PCs were extracted to represent the major axes of methylation variance across the
153 individuals and inputted into HDBSCAN (12).

154 **2.3 Stack Two – Latent age and cluster reliability**

155 Stack two produced a continuous variable called the latent age and a set of cluster
156 reliability weights to quantify the effectiveness of the HDBSCAN clusters for the prediction of
157 the age-related methylation structure. For the calculation of the latent age the chronological
158 age and epigenetic age were combined into as single latent age value using the optimised
159 weighting parameter using the 1D Fisher optimisation problem (13).

160 The final reliability metrics were averaged for production of a final reliability score for
161 each cluster. These weights were used within Stack Three for the Hierarchical Bayesian model.

162 **2.4 Stack Three – Bayesian Hierarchical Model**

163 Stack Three was determined by a Bayesian Hierarchical model for the estimation of the
164 relationship between latent age and maximum aerobic dive capacity by accounting for the
165 cluster level structure and the reliability weights generated for each sample (Stack Two). The
166 model was weighted by the cluster reliability score for the given animal within the cluster.

167 **2.5 Stack Four – CpG-level EWAS via ADVI**

168 Stack Four performed CpG EWAS with the Bayesian residuals from Stack Three as the
169 outcome of the ADVI-EWAS model (14)

170 **2.6 Stack Five – Horseshoe Bayesian shrinkage model:**

171 Stack Five computed the singular-CpG score for the global-local shrinkage model
172 (Horseshoe prior) and the estimation of methylation effects upon the maximum aerobic diving
173 capacity. The Horseshoe model was selected for its ability to retain large CpG effects while

174 aggressively shrinking noise, which is appropriate for sparse high-dimensional methylation
175 data (15).

176 **2.7 Stack Six – Overlaid Manhattan and Volcano Plots:**

177 Stack Seven incorporated the univariate CpG-specific EWAS results from Stack Four
178 and the multivariate Horseshoe shrinkage estimates from Stack Five for the generation of
179 comparative Manhattan and Volcano plots.

180 **2.8 Stack Seven – Sliding Window CpG aggregation model and gene-phenotype 181 association visualisation**

182 Stack Seven represents the primary methodological contribution of the paper. It
183 aggregated the CpGs within chromosomes in a computationally efficient manner (sliding
184 windows) by self-tuning to identify the shortest genomic segments of neighbouring CpGs
185 whose combined methylation signal showed a consistent and directionally stable association
186 with maximum aerobic diving capacity. Consistent and directionally stable in this context refers
187 to an agreed upon impact of methylation upon the aerobic diving capacity – be it hypo- or
188 hypermethylation. CpGs were ordered by chromosome and genomic positioning. For each CpG
189 the window was initiated at the CpG before being expanded iteratively to the downstream CpGs
190 provided that: (1) both the ADVI-derived effects (γ_j) and the Horseshoe-derived effects (β_j)
191 produced an increasing variance (R^2) explained by the window – with penalisation for the
192 window size – and (2) both the ADVI model and the Horseshoe model must agree on the
193 direction of association with the phenotype (methylation producing an increase or decrease in
194 the maximum aerobic diving capacity). Windows failing to meet these criteria or not reaching
195 the minimum CpG count were removed.

196 Each retained window was assigned a status: trusted, semi-trusted or weak. The status
197 was determined by the joint explanation of variance (penalised R^2 support from the ADVI and

198 Horseshoe models. Only the trusted windows were deduplicated and mapped to the genes of
199 *Z. californianus* through the GREAT-style regulatory domain definition – for the GFF3 gene
200 annotation (16). Genes were only assigned to a window when the window’s genomic span
201 overlapped the gene’s regulatory domain.

202 The transformed trusted windows were converted into gene-level and window-level
203 methylation scores and visualised for their structure across animals and clusters. Only trusted
204 windows that were unique were kept for downstream analyses to prevent the artificial inflation
205 of gene importance from duplicate gene-CpG associations. CpGs belonging to the gene-
206 assigned windows were aggregated into the per-animal scores using the logistic-scaled poly-
207 CpG transformation following the combining of the ADVI and Horseshoe effect sizes with the
208 following equation:

209

$$210 \quad \theta_j = \frac{\gamma_j + \beta_j}{2}$$

211 The scores were further assembled into two matrices; animals x genes and animals x
212 windows. The corresponding cluster-level matrices were generated with the same process but
213 for the clusters rather than animals.

214 Gene-level directional effects were generated by computing the average combined effect
215 size $(\gamma_j + \beta_j)$ across the CpGs belonging to each gene. The adaptive median-absolute-deviation
216 threshold was utilised to determine if the genes were contributing positively (Up), negatively
217 (Down) or not at all (Neutral). A positive contribution was translated as methylation increasing
218 the phenotypic trait. The per-animal gene-level effects were represented as a linear regression
219 model for the maximum aerobic dive capacity to visualise the phenotype-linked trends for the
220 selected genes. Sliding windows were used to capture regional methylation patterns, which

221 often exert biological effects through coordinated CpG behaviour rather than isolated CpG
222 changes. Further detail is provided within Supplementary Methods 1.10.

223 To evaluate the performance of the sliding window framework under varying levels of
224 noise and sparsity, a synthetic methylation dataset was generated across 3 x 3 grid of noise (0.5,
225 1.0, 2.0) and sparsity (0.2, 0.4, 0.8). Each dataset consisted of 200 CpGs with randomly
226 assigned genomic positions for 120 simulated ‘animals’. The portion of significant CpGs was
227 determined by the sparsity parameter for the assignment of the true non-zero effects of both the
228 ADVI-EWAS and the Horseshoe parameters. Methylation matrices were assumed to have a
229 standard normal distribution with phenotypes resulting from the linear combination of the ‘true’
230 (simulated) ADVI-EWAS values plus the Gaussian noise. The estimated ADVI-EWAS and
231 Horseshoe values were obtained by adding small Gaussian uncertainty to the ‘true’ effects. The
232 posterior probabilities for both statistical frameworks were computed for the probit model and
233 the CpGs were considered significant for both posteriors when exceeding 0.5 with a matching
234 posterior direction. The sliding window algorithm was performed as described above. For each
235 simulation the number of trusted windows was defined as above, the number of recovered
236 windows was defined as at least one CpG overlapping with the window being significant, the
237 precision, median window size and median R^2 was determined for every window.

238 **2.9 Stack Eight - Methylation trajectory inferencing and phenotype alignment:**

239 Stack Eight constructed a methylation-based developmental trajectory for the
240 methylation profiles to visualise the association of individuals along a continuum. The
241 continuum was aligned against the age, phenotype and gene-level methylation effects. The k-
242 NN algorithm was used to cluster the data, the MST was used to construct a relative tree for
243 the determination of association (17) and the DPT was used to provide a pseudotime-frame
244 (18).

245 **2.10 Overview of the computational processing:**

246 All Stacks were conducted in python-3.13 (Python Software Foundation, Wilmington,
247 DE, USA), on a laptop – except for Stack Four and Stack Five which were performed within
248 Google Colab’s python platform (Google LLC, Mountain View, CA, USA). The entire pipeline
249 took under two days to execute with this combination. The Hierarchical Bayesian pipeline of
250 Stack Three involved 8,000 tuning iterations and 16,000 draws. Stack Four processed 10,000
251 CpGs in batches of 1,000 CpGs for 20,000 iterations and 2,000 posterior samples. Stack Five
252 involved 30,000 iterations and 2,000 draws, thus the entire process was highly computationally
253 efficient.

254 The following python packages were used: pandas (3.0.3), numpy (2.2.6), matplotlib
255 (3.10.8), seaborn (0.13.2), scipy (1.17.1), pymc (5.28.4), pytensor (2.38.2), arviz (0.22.0),
256 statsmodels (0.14.6), hdbscan (0.8.44), scikit-learn (1.9.0) and network (3.6.1).

257 All accompanying code and data is available in the open-access GitHub repository
258 associated with this study: [tsto3616/poly-CpG: The repository encodes the sliding window
259 pipeline, accompanying statistics and data for a study on the California Sea Lion's epigenome.](https://github.com/tsto3616/poly-CpG)

260 **3. Results:**

261 **3.1 Blood methylation clusters of the *Zalophus californianus* are independent of age but 262 predictive of aerobic diving capacity:**

263 Four blood methylation profile clusters were formed, with the -1-cluster representing
264 noise (Figure 1A). The clusters did not differ significantly for chronological or epigenetic age
265 (Figure 1A; MANOVA, Roy’s greatest root, $p < 0.05$; otherwise, *ns* for all other test statistics).

266 Despite age independence the clusters were significantly differentiated for the
267 Hierarchical Bayesian and experimentally derived maximum aerobic dive capacity (Figure 1B;
268 MANOVA, $p < 0.05$).

269 **3.2 No significant CpGs explain the maximum aerobic diving capacity of *Zalophus*** 270 ***californianus*:**

271 Both univariate and multivariate CpG-phenotype scoring frameworks were applied to the
272 study of the phenotype (maximum aerobic dive capacity) – both frameworks failed to identify
273 any CpGs exceeding the significant threshold (Figure 2A, B). Effect sizes were small and
274 consistent between the two frameworks; however, the probability of the CpG impacting upon
275 the phenotype was significantly reduced for the multivariate CpG-phenotype scoring
276 framework (Figure 2A, 2B). This justified the need for the construction of a poly-CpG-
277 phenotype scoring framework.

278 **3.3 Three genes inherently linked to the maximum aerobic diving capacity and the** 279 **lactate dynamics of *Zalophus californianus*.**

280 The sliding window framework identified three windows, which were uniquely mapped
281 to the three genes: Glycine Cleavage System H (GCSH), DNA Methyltransferase 1
282 (DNMT1) and LIM Homeobox transcription factor 1- α (LMX1A). The self-tuning behaviour
283 of the sliding window framework is illustrated for the three genes and the respective genomic
284 order of the CpGs in Supplementary Figure 1.

285 All three genes were significant predictors of maximum aerobic diving capacity (Table
286 1) – consistent with either a negative or positive effect from methylation (Table 1; Figure
287 3A). The three genes each explained at least 40% of the respective variance (penalised R^2) for
288 the maximum aerobic diving capacity (Table 1).

289 The GCSH gene was independent of the DNMT1 and LMX1A genes, whereas DNMT1
290 and LMX1A both retained coordinated methylation patterns (Figure 3B).

291 **3.4 The sliding window framework is dependent on the noise: signal ratio for optimal** 292 **performance.**

293 The influence of noise and sparsity upon the performance of the sliding window
294 framework was evaluated through Python simulated datasets for a grid of noise and sparsity
295 values. At high signal levels (sparsity = 0.8), the model produced numerous trusted windows –
296 but the exact number of windows varied with noise (Figure 4A). The sparsity values of 0.4 and
297 0.2 produced zero to two windows regardless of the noise (Figure 4A).

298 The difference in the window counts of the framework maintained a consistently low
299 false-positive rate across the three noise to signal regimes (Figure 4B). The low noise-high
300 sparsity condition produced four times as many windows as the intermediate condition – both
301 conditions were determined to both present with perfect precision for the windows and an
302 acceptable rate for the inclusion of insignificant, singular-CpGs (Figure 4B; Table 2). The
303 summary statistics yielded the highest median R^2 value for the medium noise: signal regime –
304 despite similar window size (Table 2).

305 The variability in the window detection was not attributed to the differences in the
306 underlying γ or β estimates from the simulated ADVI-EWAS or Horseshoe models (Figure 4C,
307 D). R^2 distributions for γ and β parameters did not differ in a meaningful manner between the
308 two modelling approaches (Figure 4E, F). The dominant driver of window recovery was
309 sparsity. Higher sparsity produced stronger and more coherent signals for window selection,
310 whereas low sparsity reduced the ability of the framework to identify stable windows. This
311 effect was noise agnostic (Figure 4E, F; Table 2).

312 **3.5 Hierarchical analysis of the embedded methylation parameters not predictive of age**
313 **or aerobic diving capacity.**

314 Diffuse pseudotime revealed no genome-wide methylation trajectory association with the
315 phenotype, but the maximum aerobic diving capacity was strongly associated with the three
316 genes (Figure 4).

317 **4. Discussion:**

318 The aim of the paper was to develop a self-tuning window framework for the analysis
319 of the environmental-methylome-phenotype axis of 69 *Z. californianus* individuals. To do so,
320 blood samples were re-analyzed for the marine mammal specific phenotypic performance –
321 the maximum aerobic diving capacity (6, 10). The results were consistent of no phenotype-
322 methylation association when considering individual CpGs for both the univariate ADVI-
323 EWAS and the multivariate Bayesian Horseshoe shrinkage model (14, 15). Presumably the
324 results were due to the sparsity of the data, combined with the polygenic nature of the
325 complex phenotypic trait assessed and the relatively small sample size (n = 69). However, the
326 application of the sliding window framework provided vital insights into the environmental-
327 methylome-phenotype axis by identifying three key genes associated with the phenotype.
328 Thus, challenging the dominant assumption of methylome studies in mammals – the
329 methylome is more complex than predominantly being a function of age (19, 20).

330 As Sea Lions age the development from a maternally dependent pup to an independent
331 adult involves the ontogenic adaption towards the development of the physiological extremes
332 needed for prolonged diving (10, 11). However, the age of the individual is said to be the
333 primary predictor of this evolutionary dynamic. The California Sea Lion's methylome is a
334 strong predictor of aging (6, 10). A strong mechanistic explanation of these physiological
335 changes can be explained through the methylation of the three genes identified. The GCSH

336 gene is consistent with the development of post-diving hyperglycemia in pinnipeds (21-23).
337 Deficiency of the GCSH gene in humans results in non-ketogenic hyperglycemia due to the
338 accumulation of the gluconeogenic amino acid, glycine (24, 25). The excess glucose would
339 drive a metabolic shift towards lactate production in the presence of anaerobic respiration at
340 the end of the dive – resulting in potentially harmful physiological effects (26).

341 The lactate accumulation due to hypoxia may result in the phenomena known as
342 lactylation. Lactylation is a well described pathway of gene expression influenced by lactate
343 dynamics. Central to the regulation of methylation is the DNMT1 gene – lactate upregulates
344 the expression of DNMT1 due to the hypermethylation of a gene with DNMT1-suppressive
345 properties (27). DNMT1 also has protective factors against the induction of hypoxic induced
346 apoptosis (28). DNMT1 reduces the methylation rate of LMX1 – explaining the
347 interdependence of the poly-CpG-phenotype association disclosed by the sliding window
348 framework (29). LMX1A knock-out results in dysregulation of mitochondria, producing
349 excessive reactive oxygen species (ROS), which can be cytotoxic (30). Taken together, the
350 three genes are mechanistically linked to the metabolic changes experienced in response to
351 the physiological demands of diving.

352 The sparsity of the methylome reduced the power effect of the study, which was evident
353 in the lack of singular CpGs associated with the phenotype (31). However, the inclusion of
354 windows aggregated the regional CpG scores, shifting the analysis from noise abundant and
355 numerous CpGs to a coordinated regional analysis approach. The process improved the
356 reliability penalized R^2 value and directionality of the effect. As part of the model, the self-
357 tuning sliding window analysis avoids arbitrary definitions of windows by limiting the
358 windows to real chromosomes and imposing a penalization on the R^2 for the window length
359 to prevent the detection of CpGs over vast distances within the one window. The end-result is
360 phenotype-associated windows that are short, stable and directional in their predicted effect

361 upon the phenotype. Although the study was applied to one marine mammal species only, the
362 method is broadly applicable due to its ability to handle sparse data and polygenic traits (32).

363 The simulation experiments demonstrated that the sliding window framework behaved
364 predictably despite the varying noise to signal regimes. The high sparsity value produced the
365 most trusted windows, the intermediate level of sparsity generated few windows, and the low
366 level of sparsity generated no windows regardless of the noise level. For all settings
367 generating a window, there was perfect window level precision; however, false positives were
368 detected among the CpGs of windows. Despite this the framework is highly resistant to false-
369 positive calls despite substantial noise. The ADVI-EWAS and or Horseshoe model
370 parameters and their estimated R^2 distributions were not determinates of the number of
371 windows selected, nor the variability within a window. Rather sparsity was the dominant
372 determinant of window recovery – with denser signals capturing windows at a higher rate –
373 as was expected given the penalisation of the window size. The penalisation of the window is
374 inherently restricted by the density function (sparsity) of the CpG data. Future fine tuning of
375 the model could incorporate a hyperparameter for the weighting of the CpG sparsity of the
376 genome. This would be beneficial in the context of ecological studies, as the methylome
377 sequencing coverage directly impacts upon the sparsity of the methylome (33). Wildlife
378 studies are frequently limited by suboptimal sample preservation due to situations in the field
379 – a model capable of accounting for the sparsity of the CpG rate per genomic nucleotide
380 would offer a novel approach to counteract the coverage loss resulting from suboptimal
381 preservation of DNA in the field (34).

382 Despite numerous limitations throughout the study, the sliding window framework still
383 managed to produce a mechanistic explanation for the ontogenetic establishment of the
384 extreme physiological shifts towards adulthood in the California Sea Lion. The limited
385 sample size, does reduce the inferencing capabilities of the study, but it establishes the ability

386 of the sliding window framework to perform under both data and statistical sparsity (31). The
387 limitation of tissue sampling to the blood of pinnipeds is also a significant limitation.
388 Pinniped red blood cells lack nuclei and thus the epigenome is limited to cells that primarily
389 do not participate in the aerobic diving capacity – through oxygen reservation – rather the
390 effects are more likely secondary to physiological stressors rather than age (35). Examination
391 of other tissues – such as muscle – would benefit in the detection of genes associated with the
392 aerobic diving capacity of *Z. californianus*, but it would also provide stronger effect sizes
393 (10). The stronger effect sizes do not necessarily reflect the sparse nature of the data. Blood
394 samples are also one of the most accessible tissues for wildlife DNA sampling – allowing for
395 more meaningful comparisons for a pan-wildlife methylation tool (6, 34). However, blood is
396 not a universal window into whole animal methylation – tissue dependent methylation is well
397 described (36).

398 Limitations that detract from this study’s applicability include the determination of the
399 aerobic dive capacity as a function of age and the cross-sectional nature of the data recycled
400 in the study. Given that the age of *Z. californianus* is a strong predictor of the aerobic dive
401 capacity and the molecular clock of individual pinnipeds, the methylation profiles predictive
402 of the aerobic dive capacity are inherently associated with the age of the individuals (6, 10).
403 This is further exacerbated by the estimation of the maximum aerobic dive capacity of *Z.*
404 *californianus* being calculated as a linear function of the age of the individuals sampled for
405 the methylome sequencing (10). The clusters identified were age independent despite the
406 genes’ effect sizes being determined in part by cluster. Further research is a unique
407 opportunity to apply the framework to more diverse phenotypic traits that are (i) documented
408 as part of the sampling framework, (ii) are diverse within a species and (iii) are applicable
409 across species. Through future research a more mechanistic understanding of conservation
410 physiology and genetics can be established.

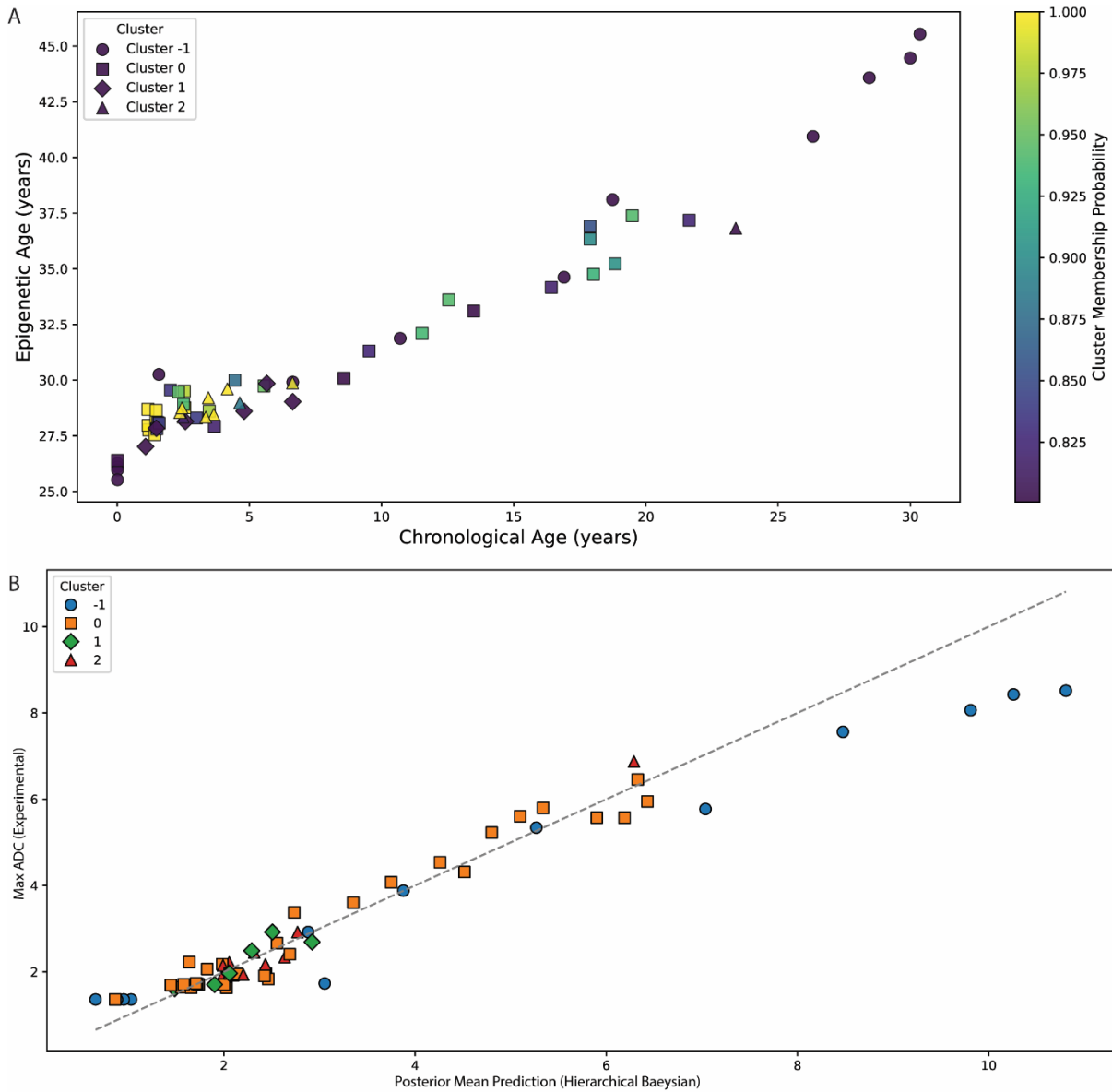
411 **5. Conclusions:**

412 The aim of the study was to develop a framework for poly-CpG regional analysis to
413 improve the understanding of methylation-phenotype associations in wildlife. Through a
414 sliding window framework, the regional poly-CpG effects were characterized for California
415 Sea Lions – three windows were isolated producing three explanatory genes. The
416 applicability of the method is widespread and can be applied to various wildlife species.

417 **6. Acknowledgements:**

418 No funding was provided. The Author would like to thank the academic advisors who
419 did not contribute to the study, but provided much needed mentorship over the years.

420



422

423 **Figure 1: Hierarchical Bayesian prediction of the maximum aerobic diving capacity**

424 **defined as a function of the experimentally derived maximum aerobic diving capacity.**

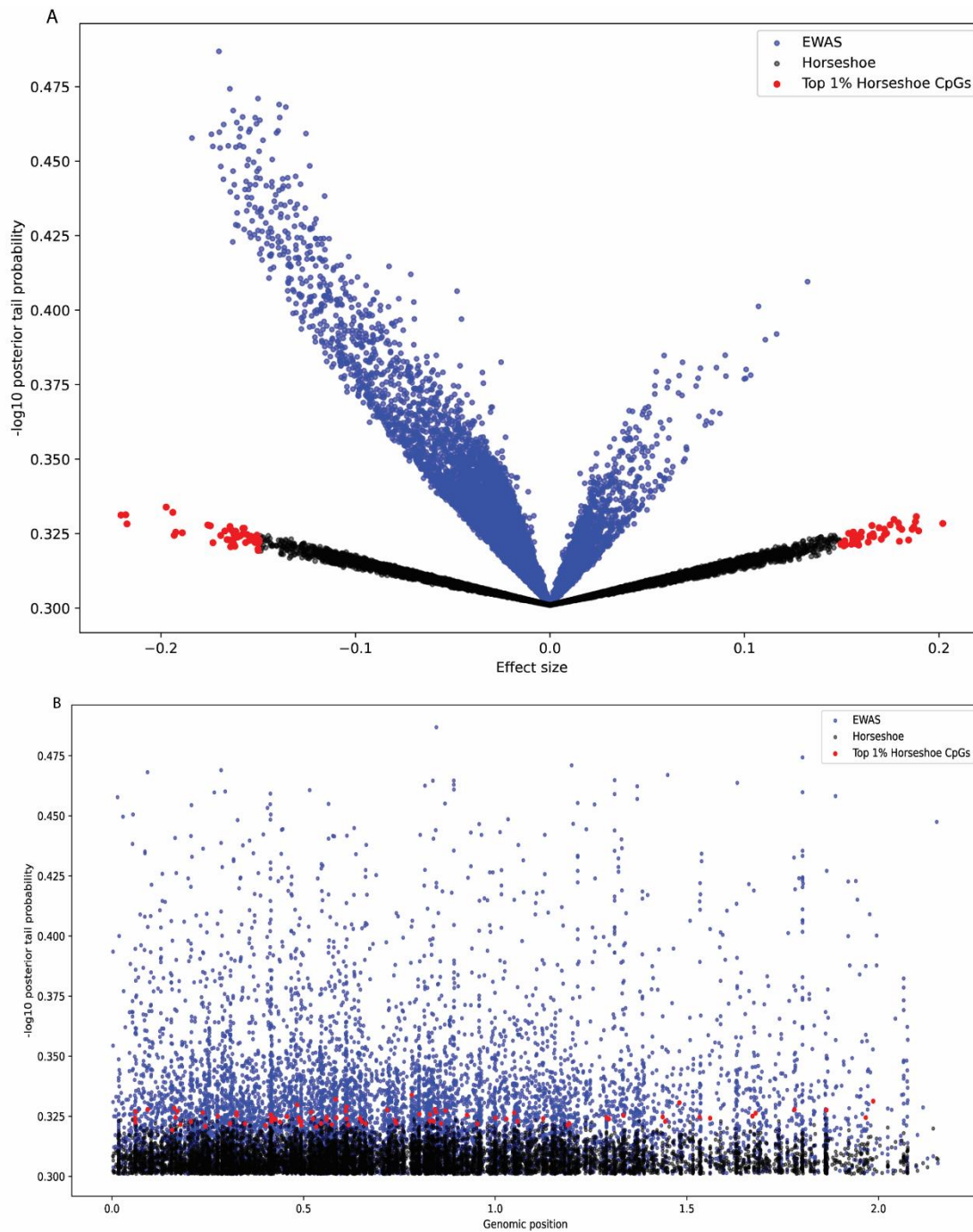
425 A. The chronological age (x axis) is compared to the epigenetic age (y axis), with the cluster

426 (shape) and probability of assignment to a cluster (colour) featured. B. The Maximum

427 experimental aerobic dive capacity (Max ADC – y axis) is defined for the mean posterior

428 Hierarchical Bayesian prediction of the ADC (x axis). The dashed line represents the perfect

429 prediction, not the linear regression line.



430

431 **Figure 2: CpG scores are not a significant predictor for the aerobic diving capacity of**

432 *Zalophus californianus*. For both A and B, the y axis features the $-\log_{10}$ posterior tail of

433 probability and the same colour coding system is retained. The black dots are the bottom 99%

434 horseshoe values, the red dots the top 1% values and the blue dots are the ADVI-EWAS values.

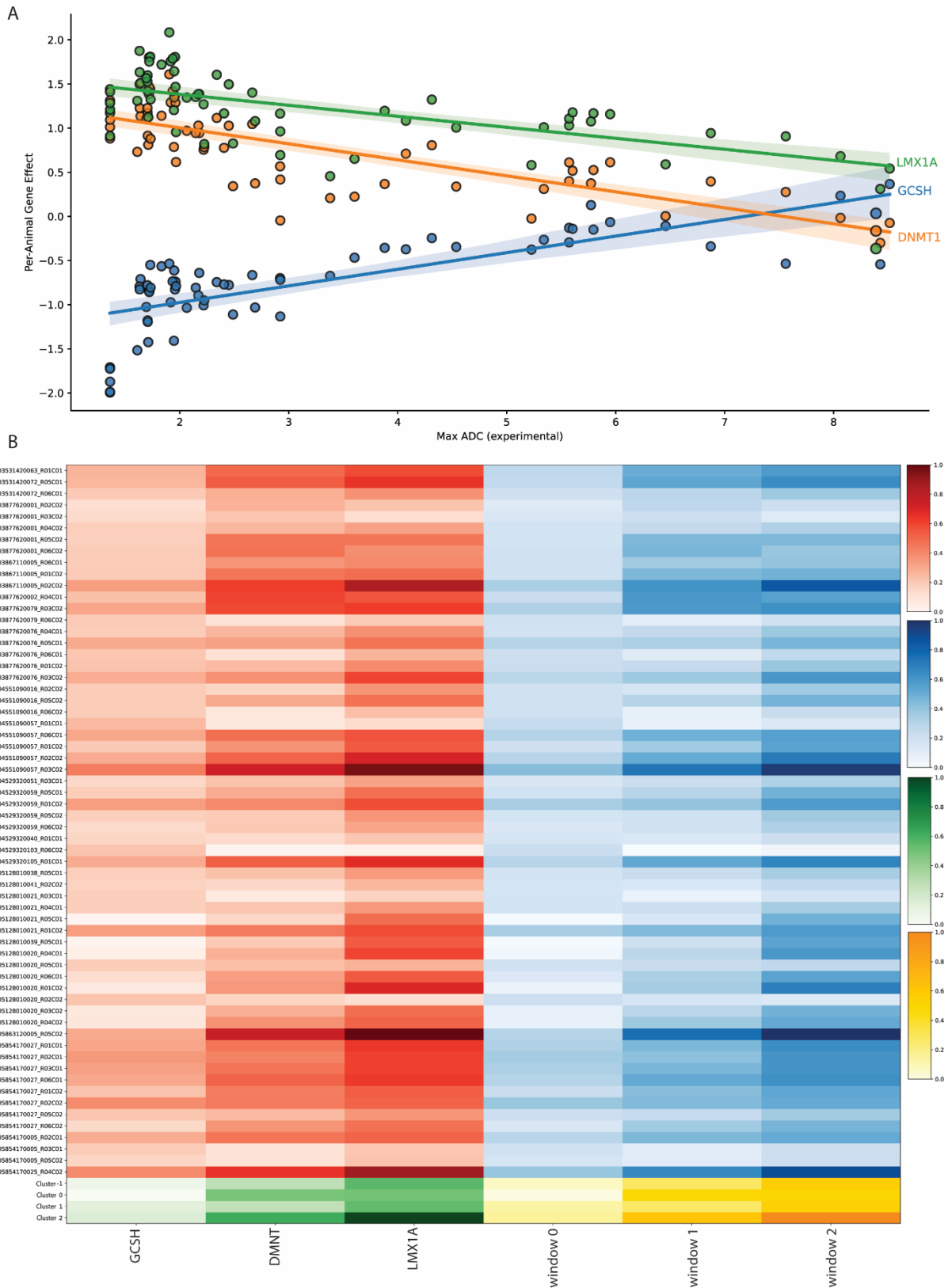
435 A. The Volcano plot of two contrasting analyses is provided for the ADVI-EWAS and the

436 Horseshoe shrinkage model. The x-axis features the effect size. B. The ADVI and Horseshoe

437 CpG values were compared within a Manhattan plot with the x axis featuring the genomic
438 position.

439 **Table 1: The GCSH, DNMT1 and LMX1A genes are all significant predictors of**
440 **maximum aerobic diving capacity.** The summary statistics for the linear regression used in
441 Figure 3A are summarised. The slope and intercept predict the per-animal, gene effect.

Statistics:	GCSH	DNMT1	LMX1A
Slope	0.1880	-0.1812	-0.1247
Intercept	-1.3511	1.3657	1.6327
R²	0.5537	0.5969	0.4148
p values	<0.0001	<0.0001	<0.0001
Standard error	0.0218	0.0192	0.0191

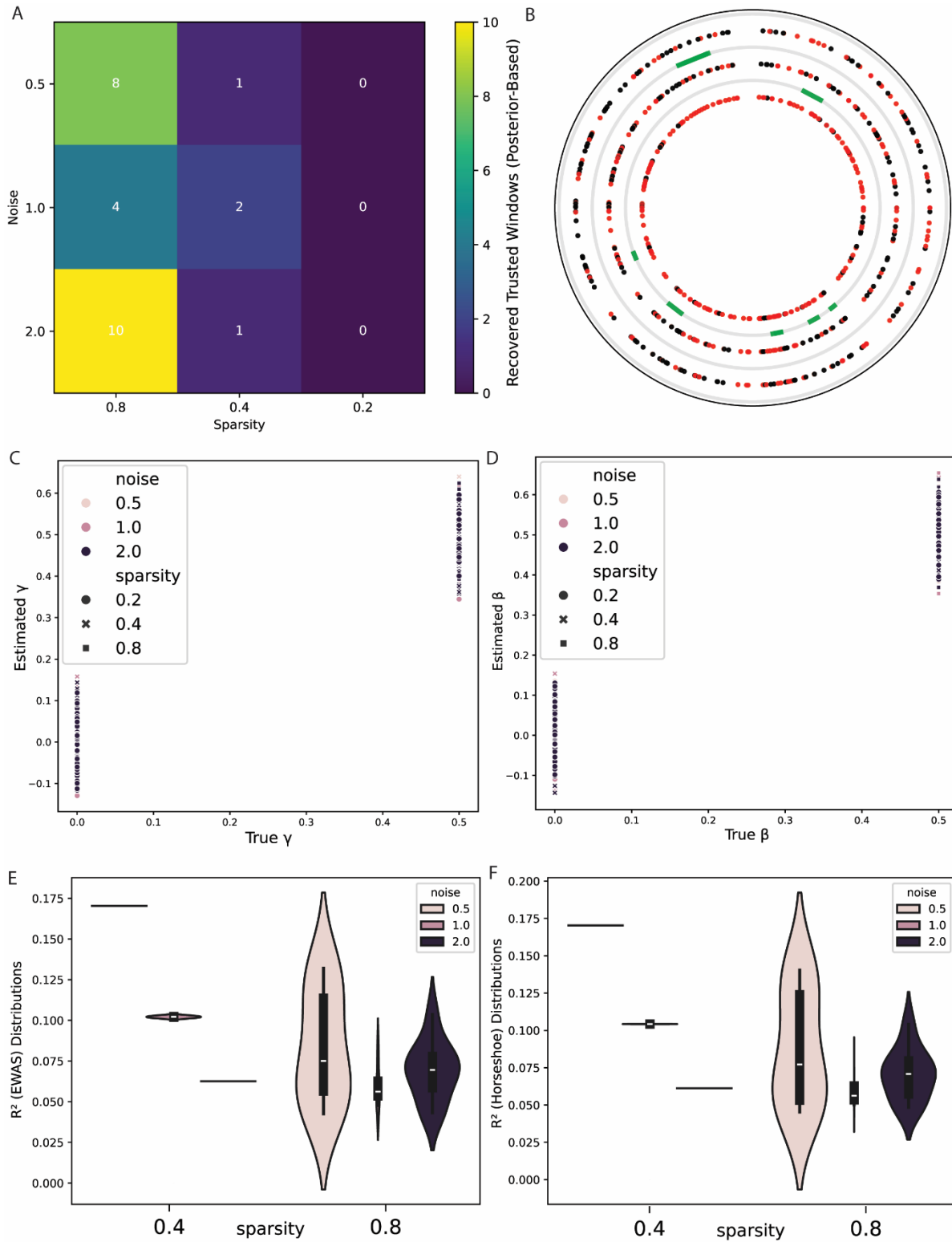


442

443 **Figure 3: Three genes strongly associated with maximum aerobic diving capacity.** A. The

444 effect sizes of the genes were mapped through linear regression for the maximum aerobic

445 capacity – Max ADC. B. A heat map was generated for the relative effect of each gene upon
446 animals (top left quad - red), window level methylation scores for the animals (top right quad
447 - blue), gene upon HDBSCAN clusters (bottom left quad - green) and window upon the
448 HDBSCAN clusters (bottom right quad - yellow).



449

450 **Figure 4: Simulations of the sliding window analysis reveal an absence of window**
 451 **identification at high noise to signal ratios and hallucinations at low noise to signal ratios.**

452 A. The heatmap shown is coloured for the number of trusted windows discovered for three tiers

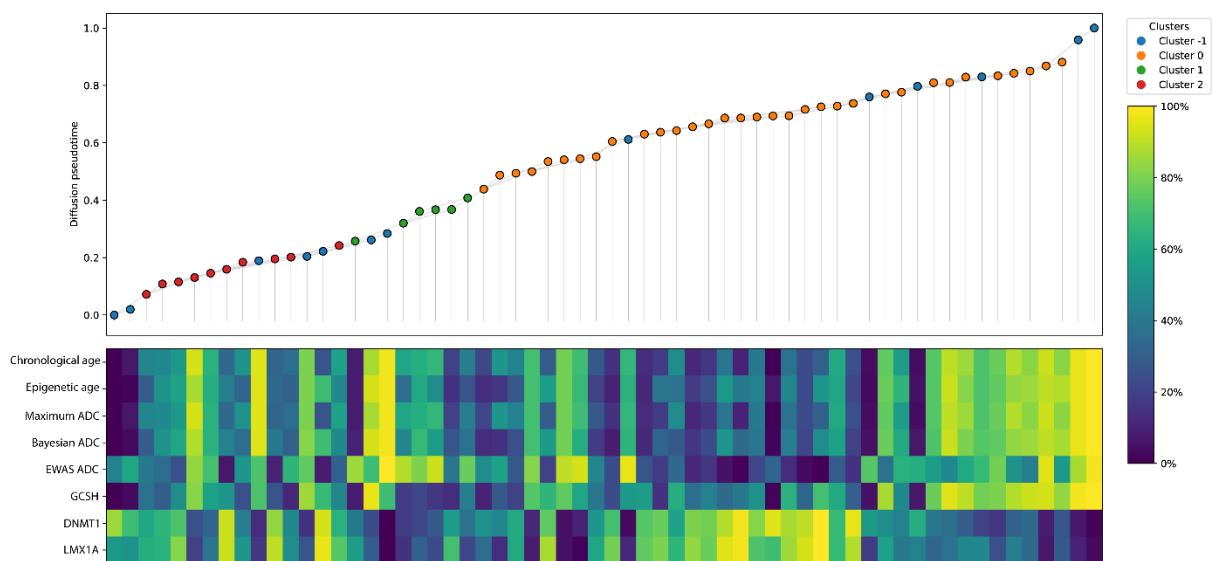
453 of sparsity (0.8, 0.4, 0.2) and noise (0.5, 1.0, 2.0) – where the lower values of sparsity were
454 predictive of a high noise: signal ratio, higher values the opposite. For the noise values the
455 higher the value the greater the noise: signal ratio. The numbers inside the tiles are the number
456 of windows identified. B. The circo plot was fitted for three combinations of the sparsity and
457 the noise were selected for a small noise: signal ratio (0.8 sparsity and 0.5 noise), medium
458 noise: signal (0.4 sparsity and 1.0 noise) and high noise: signal ratio (0.2 sparsity and 2.0 noise).
459 The inner circles of red and black dots and the green bars represent the small noise: signal ratio
460 dataset, the middle set of rings were for the medium noise: signal ratio and the outer rings of
461 dots and absence of green bars is the high noise: signal ratio (for which there were no identified
462 windows). The red dots represent the signal (CpG scores) and the black dots represent the noise
463 (CpG scores). The green bars are trusted windows. C. The γ values from the ADVI-EWAS
464 significance scoring from the simulated inputs/ ‘true’ values were modelled on the x axis and
465 the sliding window estimated γ values were plotted on the y axis. This plot was repeated for D.
466 but with the estimated β (y axis) and ‘true’ β (x axis). In plot E. the R^2 ADVI-EWAS distribution
467 (y axis) was determined for the two sparsity categories for which sliding windows were
468 successful – the 0.4 and 0.8 sparsity values (x axis). This was repeated in plot F. for the R^2
469 Horseshoe distribution (y axis) mapped against the same sparsity values (x axis).

470

471 **Table 2: The noise to signal ratio has a dramatic effect on the calling of windows and the**
472 **variance explained by the windows.** The high noise: signal ratio was for a sparsity value of
473 0.2 and a noise value of 2, the medium was 0.4 (sparsity) and 1 (noise), whereas the small
474 noise: signal ratio had a sparsity value of 0.8 and a noise value of 0.5. Detailed explanations of
475 the workflow can be found in the methods.

Statistics:	High noise: signal ratio	Medium noise: signal ratio	Small noise: signal ratio
Significant CpGs	83	101	168
Trusted windows	0	2	8
Recovered windows	0	2	8
Precision	0	1.0	1.0
Median window size	0	4.5	4.0
Median R^2	0	0.1021	0.0751

476



477

478 **Figure 5: Hierarchical ordination of the methylation does not predict phenotype.** The
479 diffusion pseudotime (from the origin – 0 to the terminal individual – 1) was mapped for a k-
480 NN, MST constructed hierarchical model for the aging and tracing of evolutionary methylation
481 histories. The origin animal at the bottom of the graph was the root for the tree, it was the
482 youngest chronological individual. The dot's colours in the top of sector of the graph represent
483 the HDBSCAN clusters. The y axis of the top plot represents the pseudotime, the x axis was
484 ordinally spaced out to accommodate for the heatmap below. The results were combined with

485 the phenotypic markers: chronological age, epigenomic age, maximum aerobic diving capacity
486 (Maximum ADC), Bayesian ADC (the predicted aerobic dive capacity from the Hierarchical
487 Bayesian model), the EWAS ADC (generated from the ADVI-EWAS model) and the GCSH,
488 DNMT1 and LMX1A genes are displayed. The colour for the heatmap dictates the relative
489 proportional value for the given trait across all samples/ individuals.

490

491 **Supplementary Methods:**

492 **1.1 Supplementary methods overview:**

493 Removal of the excess CpGs was performed using the `var()` command from Python's
494 `numpy` command. MANOVA was applied to understand the role of the chronological age and
495 epigenetic age upon the formation of the HDBSCAN clusters.

496 The 20 PCs were inputted into HDBSCAN with a minimum cluster size of six individuals
497 and minimum number of neighbouring samples of three (`min_samples = 3`) needed for the
498 identification of a core sample within a cluster. The choice of the minimum cluster size and
499 number of neighbouring samples which arbitrarily defined in this study were chosen on the
500 basis of a small sample size, high dimensional PCA space. The pipeline is described in greater
501 detail in Supplementary Methods 1.2

502 The weighted parameter (λ) was optimised using a 1D Fisher solution. Epigenetic ages
503 were calculated through using the model provided by Robeck, Haghani (6) for all pinniped
504 blood methylation clocks, as the California Sea Lion specific clock model was not freely
505 available (Supplementary Methods 1.3). The weighting parameter was chosen to maximise the
506 separation between clusters while minimised the within-cluster variability. Latent age values
507 were assigned to each individual sample and/or animal (Supplementary Methods 1.4).

508 The stability and age-related relevance of each cluster was determined through three
509 reliability metrics: (1) age coherence, (2) ordering reliability and (3) methylation tightness. Age
510 coherence was calculated as how tightly chronological and epigenetic ages were grouped
511 within a cluster. The ordering reliability was to determine if the clusters followed a monotonic
512 progression in which the epigenetic age either always increased or decreased with the ordering
513 of the hierarchical clusters. The methylation tightness quantified the within-cluster

514 homogeneity of the CpG M-values. The Hierarchical Bayesian model produced a divergence
515 rate of 0.075%. This was summarised in Supplementary Methods 1.5.

516 The maximum aerobic dive capacity was assumed to be a normally distributed outcome
517 variable with a mean provided as a linear prediction generated from an intercept, the latent age
518 effect and a random effect to account for HDBSCAN clustering. Individual weights were
519 included in the model as a scale for the impact of the cluster reliability for the given animal.
520 That is the more 'reliable' the cluster the animal belongs to, the more the animal contributes to
521 the scaling of the residual variance and thus contribute more strongly to the likelihood
522 prediction for the Bayesian model. For all informative priors, weak parameters were used to
523 prevent overfitting without over-constraining the model – a necessary consideration or the
524 limited sample size of the study (n=69). Posterior samples were estimated using Markov chain
525 Monte Carlo (MCMC). The posterior mean predictions were used to compute the residuals
526 needed for the downstream ADVI modelled EWAS in Stack Four. The formulae for the
527 Bayesian Hierarchical model are available in Supplementary Methods 1.6, with additional
528 information regarding the MCMC sampling.

529

530 The residual maximum aerobic dive capacity was determined as a function of latent age,
531 CpG-specific methylation, cluster-level random effects using the individual reliability weights
532 from Stack Two. To computationally manage the CpG EWAS (10,000 CpGs, 20,000 iterations
533 and 2,000 posterior samples) the process was conducted in batches of 1,000 CpGs with separate
534 regression fit for each CpG using ADVI rather than a MCMC posterior sampling. The ADVI
535 process provided a faster approximation of the posterior parameters. For each of the CpGs the
536 estimated CpG-specific effect size (γ) determines the effect of the methylation upon the residual
537 phenotype. The residual phenotype is a condition of the latent age and clustering associated

538 phenotype from Stack Three. The posterior means and standard deviations for the CpG-specific
539 effect size were extracted and retained for all CpGs. These were used to compute the 95%
540 credible intervals. These were used for the determination of the posterior probabilities of a
541 positive benefit under the Gaussian variational posterior. The CpGs were classified as
542 significant for the posterior probability of an effect being greater than 0.975 for either direction.
543 The resultant CpG effect size data frame was used for downstream analyses. The statistical
544 basis of Stack Four can be found in Supplementary Methods 1.7.

545 To achieve this the methylation matrix was standardised for the columns to ensure the
546 effect sizes across the CpGs are comparable to the ADVI results of Stack Four. A Bayesian
547 linear regression model with a horseshoe prior was fit to the standardised methylation matrix
548 by applying ADVI for the computational efficiency.

549 The estimated model estimated a global shrinkage parameter (τ), singular-CpG-specific
550 local shrinkage parameters (λ_j), and singular-CpG-specific regression coefficients (β_j). The
551 Horseshoe prior enables strong reduction of the impact of noise (CpGs) while allowing for
552 large CpG effects to remain unshrunk. Posterior samples were drawn from the fitted variational
553 distribution with the posterior means and standard derivations of β_j were extracted every
554 individual CpG. Additional resources for Stack Five can be found in Supplementary Methods
555 1.8.

556 Each CpG was aligned for its EWAS posterior mean effect size ($\hat{\gamma}_j$) and its comparable
557 Bayesian tail probability for the Horseshoe model ($\hat{\beta}_j$) and posterior standard deviation. The
558 two datasets were annotated on the CpG-level with the computation of the cumulative genomic
559 positions for Manhattan-style plotting.

560 Both plots were developed by plotting the ADVI posterior, Bayesian tail probability for
561 γ_j and the posterior tail of probability for the Horseshoe model (β_j). The CpGs from the

562 Horseshoe CpGs were plotted using the assumption that the posterior tail of probability (β_j)
563 generated by the Horseshoe CpGs under variational approximation was comparable to the
564 ADVI posterior Bayesian tail probability for γ_j . The Volcano plot visualised the direction of the
565 significance for methylation while the Manhattan plot visualised the overall significance of
566 methylation for the overlaid genomic positions. The top 1% most significant Horseshoe CpGs
567 were highlighted in both plots. Comprehensive detail on the statistical foundations of these
568 steps can be found in Supplementary Methods 1.9.

569 CpG methylation values were reduced using PCA and a k-nearest neighbour (k-NN)
570 graph was constructed in PCA space with a minimum spanning tree (MST) extracted for the
571 backbone of the methylation similarity score across individuals. Diffusion pseudotime was
572 computed for the y axis of the MST and rooted for the individual with the lowest chronological
573 age. The production was a one-dimensional ordering of individuals that reflects the progressive
574 changes in the methylation structure. The MST was visualised within the PCA space for
575 clustering-based colouring (from Stack One) to highlight the relationship between the
576 methylation-defined clusters and inferred trajectory.

577 For the integration of the phenotypic information the pseudotime ordering was used as a
578 shared x-axis for the second panel. The second panel featured a tile-based heatmap of the
579 ranked phenotypes with variables such as chronological age, epigenetic age, maximum aerobic
580 diving capacity, the Bayesian residuals from Stack three, the ADVI derived poly-CpG scores
581 from Stack Four result and the gene-level effects for the genes determined to be significant for
582 the methylation-phenotype association. Each variable is percentile normalised and displayed
583 as a horizontal tile row with precise alignment against the MST ordering. The results produced
584 a unified figure for the comparisons of methylation trajectory structure against the
585 physiological and gene-level phenotypes. Trajectory inference was used to identify latent
586 methylation progression patterns that may reflect physiological ageing or adaptive diving

587 phenotypes. Further information detailing the Stack Eight pipeline can be found in
588 Supplementary Methods 1.11.

589 **1.2 Calculation of the M-scores:**

590 The β -scores from the normalised methylation data were converted into M-scores for
591 methylation using the formula:

592

$$594 \quad M_{ij} = \log_2 \left(\frac{\beta_{ij}}{1 - \beta_{ij}} \right)$$

593 **1.3 Estimation of the Epigenetic age:**

595 The epigenetic age of the animals was predicted using a pinniped-wide methylation
596 clock. For the calculation of the methylation age an elastic net epigenetic clock model was used
597 whereby the methylation age ($\text{Age}_i^{\text{epi}}$) is a function of the intercept β_0 and the sum of the
598 coefficients for the CpG (β_j) and the methylation scores (M-scores). The epigenetic age was
599 calculated using the formula:

$$600 \quad \text{Age}_i^{\text{epi}} = \beta_0 + \beta_1 M_{i1} + \beta_2 M_{i2} + \dots + \beta_p M_{ip}$$

601 **1.4 Latent age definition:**

602 The latent age is then defined by the combination of the epigenetic age, the chronological
603 age and age of the clusters from stack 1. The latent age was classified using the formula:

$$604 \quad \text{LatentAge}_i(\lambda) = \lambda \text{Age}_i^{\text{chron}} + (1 - \lambda) \widehat{\text{Age}}_i^{\text{epi}}, \quad \lambda \in [0,1]$$

605 Where the λ is determined by the 1D Fisher optimisation formula:

606

$$607 \quad \lambda^* = \arg \max_{\lambda \in [0,1]} J(\lambda)$$

608 $J(\lambda)$ is given by:

609

610
$$J(\lambda) = \frac{S_B(\lambda)}{S_W(\lambda)}$$

611 Where $S_B(\lambda)$ and $S_W(\lambda)$ – the between cluster and within cluster parameters respectively
612 – were given by:

613

614
$$S_B(\lambda) = \sum_k n_k \left(\overline{\text{LatentAge}}_k(\lambda) - \overline{\text{LatentAge}}(\lambda) \right)^2$$

615
$$S_W(\lambda) = \sum_k \sum_{i \in k} \left(\text{LatentAge}_i(\lambda) - \overline{\text{LatentAge}}_k(\lambda) \right)^2$$

616 The individual animal mean latent age was provided by the formula:

617

618
$$\overline{\text{LatentAge}}(\lambda) = \frac{1}{N} \sum_{i=1}^N \text{LatentAge}_i(\lambda)$$

619 The cluster mean latent age was provided by:

620
$$\overline{\text{LatentAge}}_k(\lambda) = \frac{1}{n_k} \sum_{i \in k} \text{LatentAge}_i(\lambda)$$

621 **1.5 Calculation of the reliability scores and weights:**

622 The formula for this reliability weighting is provided by formula below whereby the age
623 coherence (r_k^{coh}) is a function of the variance of the chronological age and epigenetic age. Age
624 coherence scores close to 1 indicate a tight age clustering for both the chronological and
625 epigenetic age, but a score of 0 indicates the opposite.

626

627

$$r_k^{\text{coh}} = \frac{1}{1 + \text{Var}(\text{Age}_i^{\text{chron}}; i \in k) + \text{Var}(\widehat{\text{Age}}_i^{\text{epi}}; i \in k)}.$$

628

The order of the ages for the clusters is defined in two steps, first by the function:

629

630

$$\overline{\text{Age}}_k = \frac{1}{n_k} \sum_{i \in k} \widehat{\text{Age}}_i^{\text{epi}}$$

631

Then by the formula produces the reliability order score:

632

633

$$r_k^{\text{ord}} = \mathbb{I}(\overline{\text{Age}}_k < \overline{\text{Age}}_{k+1})$$

634

The internal methylation tightness within the cluster is given by the formula:

635

636

$$r_k^{\text{tight}} = \frac{1}{1 + \text{Var}(M_{ij}; i \in k, j \in \mathcal{C})}$$

637

Where the r_k^{tight} value is a function of the variation of M-score for the clusters for all

638

10,000 CpGs (where \mathcal{C} is the CpGs) used for the clustering. If the clusters look similar in

639

methylation space the value is close to 1 but if not, it is close to 0. It is a molecular coherence

640

metric.

641

All the above cluster evaluation metrics are combined into a combined r_k for the

642

calculation of the weights for the Bayesian analysis.

643

644

$$r_k = \frac{r_k^{\text{coh}} + r_k^{\text{ord}} + r_k^{\text{tight}}}{3}$$

645 The reliability of the cluster (determined from r_k) is fed into the formula for the reliability
646 of each individual animal. The reliability of the cluster and the predicted likelihood of the
647 animal belonging to the cluster are combined to produce the metric:

648

649
$$w_i = p_i^{\text{cluster}} r_{c_i}$$

650 The w_i value is used as a Bayesian weight/ prior.

651 **1.6 Determination of the Bayesian Hierarchical model:**

652 The Bayesian Hierarchical model will be constructed using the formula:

653

654
$$Y_i | \mu_i, \sigma^2, w_i \sim \mathcal{N}\left(\mu_i, \frac{\sigma^2}{w_i}\right)$$

655 Where i denotes the animal, Y_i is the maximum aerobic diving capacity, the μ_i is the linear
656 predictor defined below, the σ is residual normal variation, N indicates the normal distribution
657 and w_i is the clipped weight (clipped at 10^{-6} to avoid zero variance issues) provided by the stack
658 2.

659 The σ value is defined as HalfNormal(5) whereby the residual noise is positive and
660 moderate but allows noisy data to feature. This stabilises the model with weighted likelihoods.

661 The HalfNormal(5) equation dictates that the normal distribution has a prior ranging from
662 0 to ~ 10 with most values ranging from 0-5, values above ten are possible. It serves as a weakly
663 informative prior for the prevention of unreasonable values in the absence of underfitting.

664 The formula for μ_i is provided below:

665

666
$$\mu_i = \alpha + \beta_{\text{age}} \text{LatentAge}_i + b_{c_i}$$

667 LatentAge_{*i*} is output of the optimisation from stack 2. The α value is a weak global
668 intercept prior that barely constrains the model. This prevents the drifting of the intercept due
669 to the small sample size (n=69) and allows for an effect from the data. The β_{age} latent age prior
670 is for the stabilisation of the regression model without biasing the magnitude of the effect or
671 the direction of the age effect. The b_{c_i} value is the cluster random effects defined by the cluster
672 specific deviation of the cluster c_i for the animal i .

673 The α value is defined as $\alpha \sim N(0, 10^2)$, the β_{age} latent age as $\beta_{\text{age}} \sim N(0, 10^2)$, and the b_{c_i}
674 as $b_{c_i}|\tau^2 \sim N(0, \tau^2)$.

675 The τ value is provided by the formula: $\tau \sim \text{HalfNormal}(5)$ which allows for model cluster
676 differences but also allows for extreme values if needed. It is a shrinkage mechanism to prevent
677 the overinterpretation and or overfitting from clusters and or small sample size.

678 The posterior sampling was performed using No-U-Turn Sampler (NUTS) for which
679 there were 4 chains, 2,000 tuning steps, 4,000 posterior draws, a target of acceptance of 0.9 and
680 a fixed random seed for reproducibility.

681 The posterior means were used to calculate the residuals using the following formulae:

682

$$683 \hat{\mu}_i = \mathbb{E}[\alpha] + \mathbb{E}[\beta_{\text{age}}] \cdot \text{LatentAge}_i + \mathbb{E}[b_{c_i}]$$

684 Where the residuals are:

685

$$686 \text{Residual}_i = Y_i - \hat{\mu}_i$$

687 The residuals were included in the ADVI component of the EWAS for Stack Four.

688 **1.7 Single-CpG EWAS model overview:**

689 The ADVI process involved the optimising the parameters through the maximisation of
 690 the evidence lower bound (ELBO). The posterior summaries of the predicted effect of
 691 methylation upon the aerobic dive capacity residuals (γ_j) from Stack Three were generated by
 692 the sampling of the optimised variational distribution.

693 The model assumes a global normal and linear formula for the prediction of the residuals
 694 from stack 3 using the following equation:

$$696 \quad r_{ij} = \alpha_j + \beta_{age,j} \cdot \text{LatentAge}_i + \gamma_j \cdot M_{ij} + b_{cij} + \varepsilon_{ij}, \quad \varepsilon_{ij} \sim \mathcal{N}\left(0, \frac{\sigma_j^2}{w_i}\right)$$

697 The weight is derived from the stack 2, the $b_{cij} \sim N(0, \tau_b)$ for which $\tau_b \sim \text{HalfNormal}(5)$.
 698 The linear model has the following priors for all subscript j , each CpG:

699

$$700 \quad \alpha_j \sim \mathcal{N}(0, 10^2), \quad \beta_{age,j} \sim \mathcal{N}(0, 10^2), \quad \gamma_j \sim \mathcal{N}(\mu_\gamma, \tau_\gamma^2)$$

701 The hyperpriors are defined as below:

702

$$703 \quad \sigma_j \sim \text{HalfNormal}(5), \quad \mu_\gamma \sim \mathcal{N}(0, 10), \quad \tau_\gamma \sim \text{HalfNormal}(5)$$

704 The ADVI method to estimate the γ_j values for the given age and M score of each CpG
 705 per animal (M_{ij}). The γ_j values are given by the Gaussian factorised approximation formula,
 706 where d_j is the number of parameters:

707

$$708 \quad q(\theta_j) = \prod_{k=1}^{d_j} \mathcal{N}(\theta_{jk} \mid \mu_{jk}, s_{jk}^2)$$

709 Whereby the θ_j is defined as the following:

710

712

$$\theta_j = (\alpha_j, \beta_{\text{age},j}, \gamma_j, \sigma_j, b_{1,j}, b_{2,j}, \dots, b_{K,j})$$

711

713 The variation parameters (μ_{jk} and s_{jk}) are optimised using the ELBO method – as
714 described:

715

716

$$\text{ELBO}(\phi_j) = \mathbb{E}_{q_{\phi_j}} \left[\log p(r | \theta_j) + \log p(\theta_j) - \log q_{\phi_j}(\theta_j) \right]$$

717

718

The summaries of the posteriors for the γ_j are obtained by sampling the variation approximations with the following equation:

719

720

$$\gamma_j^{(s)} \sim q(\gamma_j | \phi_j), \quad s = 1, \dots, S$$

721

722

723

724

725

726

The conversion of the Bayesian ADVI posterior values for each singular-CpG was used to calculate the mean and standard deviation of the predictions. These two metrics were fed into a Bayesian analogue of the p value to predict the significance of each CpG upon the aerobic dive capacity. This assumed a Gaussian variational posterior and produces a directional ‘p value’. Significance was defined as $(p > 0.975) | (p < 0.025)$ – for the 95% confidence intervals as defined as:

727

728

$$\gamma_{j,\text{lower}} = \hat{\gamma}_j - 1.96 \cdot \text{sd}(\gamma_j), \quad \gamma_{j,\text{upper}} = \hat{\gamma}_j + 1.96 \cdot \text{sd}(\gamma_j)$$

729

The equation for the conversion to the p values is provided below:

730

731
$$p(\gamma_j > 0) = 1 - \Phi\left(\frac{0 - \hat{\gamma}_j}{\text{sd}(\gamma_j)}\right)$$

732 Where Φ is the cumulative distribution means μ and standard deviation σ .

733 **1.8 Horseshoe Bayesian shrinkage model:**

734 The model relied on the standardisation of the M-score matrix to stabilise the shrinkage
735 behaviour of the horseshoe prior. This was performed through:

736

737
$$M_{ij}^* = \frac{M_{ij} - \bar{M}_{.j}}{\text{sd}(M_{.j})}$$

738 The Bayesian linear regression model with global-local shrinkage prior for the estimation
739 of the singular-CpG-specific methylation effects on the observed maximum aerobic diving
740 capacity is given for the individual i , and the single-CpG j :

741

742
$$y_i = \alpha + \sum_{j=1}^J \beta_j M_{ij}^* + \varepsilon_i, \quad \varepsilon_i \sim \mathcal{N}(0, \sigma^2)$$

743 With the shrinkage structure:

744

745
$$\beta_j \sim \mathcal{N}(0, \tau \lambda_j)$$

746 Where:

747

748
$$\tau \sim \text{HalfCauchy}(1), \quad \lambda_j \sim \text{HalfCauchy}(1)$$

749 With the additional priors:

750

751
$$\alpha \sim \mathcal{N}(0, 10^2), \quad \sigma \sim \text{HalfNormal}(5)$$

752 Posterior inferencing used mean-field ADVI with the variational posterior assumed to be
753 of a fully factorised Gaussian form:

754

755
$$q(\theta) = \prod_{k=1}^d \mathcal{N}(\theta_k \mid \mu_k, s_k^2)$$

756 Where the θ is:

757

758
$$\theta = (\alpha, \beta_1, \dots, \beta_J, \sigma)$$

759 The posterior summaries were generated from the 30,000 iterations of ADVI and 2,000
760 draws to form the optimised variational distributions. Posterior summaries were generated for
761 the estimate of β_j which was determined through the formula:

762

763
$$\hat{\beta}_j = \mathbb{E}[\beta_j \mid \text{data}], \quad \text{sd}(\beta_j) = \sqrt{\text{Var}(\beta_j \mid \text{data})}, \quad |\hat{\beta}_j| = |\hat{\beta}_j|$$

764 **1.9 Manhattan and Volcano plotting:**

765 For each CpG the EWAS directional Bayesian tail of probability was computed under the
766 Gaussian variational approximation:

767

768
$$p_{\gamma,j} = P(\gamma_j > 0 \mid \text{data}) = 1 - \Phi\left(0; \hat{\gamma}_j, \text{sd}(\gamma_j)\right)$$

769 Where the $\hat{\gamma}_j$ and $sd(\gamma_j)$ were denote the posterior mean and standard deviation
770 respectively. The Volcano statistics were generated through:

771

$$772 \text{volcano}_{\gamma,j} = \begin{cases} -\log_{10}(1 - p_{\gamma,j}), & \hat{\gamma}_j \geq 0, \\ -\log_{10}(p_{\gamma,j}), & \hat{\gamma}_j < 0 \end{cases}$$

773 Genomic coordinates of the CpGs were cumulatively considered for the Manhattan plot
774 as:

775

$$776 \text{position}_{ij} = \text{Start}_{ij} + \min(\text{Start}_{.j})$$

777 For the Horseshoe inputs the posterior probability of a positive effect are generated from
778 the Gaussian variational posterior – computed as:

779

$$780 p_{\beta,j} = P(\beta_j > 0 \mid \text{data}) = 1 - \Phi\left(0; \hat{\beta}_j, \text{sd}(\beta_j)\right)$$

781 The Horseshoe volcano plots were generated through:

782

$$783 \text{volcano}_{\beta,j} = \begin{cases} -\log_{10}(1 - p_{\beta,j}), & \hat{\beta}_j \geq 0, \\ -\log_{10}(p_{\beta,j}), & \hat{\beta}_j < 0 \end{cases}$$

784 The top 1% of the Horseshoe CpGs were isolated as a function of the magnitude of the
785 posterior mean:

786

$$787 |\hat{\beta}_j| > \text{percentile}_{99}(|\hat{\beta}|)$$

788 The Horseshoe genomic coordinates for the Manhattan plot were as described above for
789 ADVI.

790 **1.10 Sliding window and gene-phenotype visualisation analysis:**

791 The CpG position was determined computationally using the following notation to align
792 the CpGs – this was the basis of the sliding window model – so as to not confound biologically
793 independent data from different strands of DNA.

794

$$795 \text{cpg(pos)}_j = \begin{cases} \text{start}_j, & \text{strand} = +, \\ \text{end}_j, & \text{strand} = - \end{cases}$$

796 The methylation matrix was constructed with the CpGs sorted by the Chromosome and
797 genomic position. The CpGs, j , were computed within a window, w , through the following
798 formula to determine the ‘signal’, where the M denotes the methylation matrix, the γ_w and the
799 β_w denote the ADVI and Horseshoe values for each window:

800

$$801 S_\gamma = M_{\cdot,w} \gamma_w, \quad S_\beta = M_{\cdot,w} \beta_w.$$

802 The signal was used to calculate the portion of the variability that was explained by the
803 model. The formula for this process was:

804

$$805 R^2 = \left(\frac{\langle S - \bar{S}, y - \bar{y} \rangle}{\|S - \bar{S}\| \|y - \bar{y}\|} \right)^2, \quad \text{sign} = \text{sgn}(\langle S - \bar{S}, y - \bar{y} \rangle)$$

806 A penalised score of the R^2 was generated through the formula:

807

$$808 \text{score} = R^2 - \lambda_{\text{penalty}} \cdot |W|$$

809 This score was iteratively calculated for the formula $score_{new} - score_{best} > minimum-R^2-$
810 *increase*. In addition, both models must agree on the direction of the score – as determined by:

811

$$812 \quad \text{sign}_\gamma = \text{sign}_\beta \neq 0$$

813 If either the R^2 growth condition or the direction of the score condition failed then the
814 growth of the window ceased. Windows smaller than the predefined minimum number of CpGs
815 (four – can be altered) then the window was not retained. For a window to receive a trusted
816 status then the R^2 value for both the ADVI and Horseshoe sliding window models must be
817 greater than 0.5. For semi-trusted window status, the R^2 value must be greater than 0.5 for
818 either the ADVI or the Horseshoe sliding window model.

819 The windows were then assigned to genes using the transcription start site (TSS) – as
820 defined by the strand. The basal domain was assigned as the TSS – 5,000bp to the TSS +
821 1,000bp. Domains were extended to the midpoints between adjacent genes – with the closest
822 gene being assigned the CpG for intergenic regions. A window with the genomic span of s_w
823 and e_w (start and end of the window) were assigned to the gene, g , only when:

824

$$825 \quad \text{regulatory}(\text{start})_g \leq e_w \quad \text{and} \quad \text{regulatory}(\text{end})_g \geq s_w$$

826 The combined effect of the CpG was used for downstream calculation of the gene and or
827 window effects on the phenotype of *Z. californianus*. For the calculation fo the combined
828 effects the formula below was used:

829

$$830 \quad \theta_j = \frac{\gamma_j + \beta_j}{2}$$

831 The logistic transformation was applied to each animal's methylation measurement:

832

833
$$\text{poly}_{ij} = \sigma(\theta_j M_{ij}), \quad \sigma(x) = \frac{1}{1 + e^{-x}}$$

834 The combined effect per gene was defined as:

835

836
$$\text{Effect}_g = \frac{1}{|W_g|} \sum_{j \in W_g} (\gamma_j + \beta_j)$$

837 The median absolute deviation (MAD) was used to classify the direction of impact – be

838 it negative, positive or absent. For this calculation the following formulae were used:

839

840
$$\text{Effect}_g = \frac{1}{|W_g|} \sum_{j \in W_g} (\gamma_j + \beta_j), \quad \text{MAD} = \text{median} (|\text{Effect}_g - \text{median}(\text{Effect})|)$$

841 For which the direction is determined by Effect_g being greater than $0.5(\text{MAD})$ producing

842 an Up direction, an Effect_g less than $0.5(\text{MAD})$ producing a Down direction, otherwise

843 producing a Neutral direction.

844 The linear regression model for the effect of the genes upon the maximum aerobic diving

845 capacity (ADC) is predicted by the following equation:

846

847
$$E_{ig} = \alpha_g + \beta_g \text{ADC}_i + \varepsilon_{ig}$$

848 For the slope and intercept values:

849

850
$$(\hat{\alpha}_g, \hat{\beta}_g) = \arg \min_{\alpha, \beta} \sum_{i=1}^n (E_{ig} - \alpha - \beta \text{ADC}_i)^2.$$

851 The E_{ig} value in the global linear model represents the residual error term. The regression
 852 model was fitted for each individual gene.

853 The simulation for the sliding windows involved the mapping of artificial CpGs to a
 854 100,000 nucleotide-long chromosome (the ‘true’ effects). For each sparsity level (0.2, 0.4, 0.8)
 855 the proportion of the CpGs determined to be significant was given by the formula:

856

857
$$\gamma_{\text{true},i} = \beta_{\text{true},i} = \begin{cases} 0.5 & \text{if CpG } i \text{ is selected as causal,} \\ 0 & \text{otherwise} \end{cases}$$

858 Where the γ and β values are substitutes for the ADVI-EWAS and Horseshoe parameters
 859 respectively. Methylation values were drawn from 120 individual simulated samples for which
 860 there was a presumption of a standard normal distribution:
 861 $M_{n \times p} \sim \mathcal{N}(0,1)$. Phenotypes were generated as a linear combination of the ADVI-EWAS
 862 effect with the inclusion of Gaussian noise – as defined by:

863

864
$$y = M\gamma_{\text{true}} + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2)$$

865 Where σ was 0.5, 1.0, 2.0 – the values of the noise quantification. The estimation of the
 866 parameters was carried out through the adding of Gaussian uncertainty to the ‘true’ effects:

867

868
$$\gamma_{\text{est}} = \gamma_{\text{true}} + \mathcal{N}(0,0.05), \quad \beta_{\text{est}} = \beta_{\text{true}} + \mathcal{N}(0,0.05)$$

869 The posterior probabilities were determined with the following equations:

870

871
$$P(\gamma > 0) = \Phi\left(\frac{\gamma_{est}}{0.05}\right), \quad P(\beta > 0) = \Phi\left(\frac{\beta_{est}}{0.05}\right)$$

872 Where Φ is the standard normal cumulative distribution function. CpGs were considered
 873 significant when both posterior probabilities were greater than 0.5 and both signs of the γ_{est} and
 874 β_{est} agreed. The sliding window approach was performed as described above, only
 875 without consideration of the chromosome since only one ‘chromosome’ was processed.

876 **1.11 Trajectory based methylation profiling and phenotype association:**

877 Principal component analysis was applied to the methylation matrix, M:

878

879
$$M_{PCA} = PCA_{20}(M)$$

880 Only the first two components were analysed further. A k-NN graph was constructed
 881 using Euclidean distances for the methylome. For each sample, i , the edges were added to its
 882 k=10 nearest neighbours and the minimum spanning tree (MST) generated. This provided a
 883 sparse backbone for the analysis of the methylome across individuals:

884

885
$$G_{MST} = MST(G_{k-NN})$$

886 The root node was chosen as:

887

888
$$i_{root} = \arg \min_i Age_i^{chrono}$$

889 Where the Age_i^{chrono} is the chronological age. The diffusion pseudotime was computed
 890 along the shortest path distance of the MST:

891

892
$$\text{DPT}_i = d(i_{\text{root}}, i)$$

893 The results were normalised to [0, 1]:

894

895
$$\text{DPT}_i^{\text{norm}} = \frac{\text{DPT}_i - \min_j \text{DPT}_j}{\max_j \text{DPT}_j - \min_j \text{DPT}_j}$$

896 The PCA-MST embedding was generated by the plotting of the animals in the PC1, PC2
897 plane. The edges (u, v) were drawn as line segments and the nodes coloured by cluster labels
898 from Stack One. Individuals were ordered according to increasing methylome pseudotime with
899 the provided formula:

900

901
$$\pi = \text{argsort}(\text{DPT}^{\text{norm}})$$

902 For each phenotype the values were rank normalised as described below:

903

904
$$\tilde{v}_i = 100 \times \frac{\text{rank}(v_i) - 1}{n - 1}$$

905 The heatmap was generated with rows corresponding to the individuals ordered by
906 pseudotime:

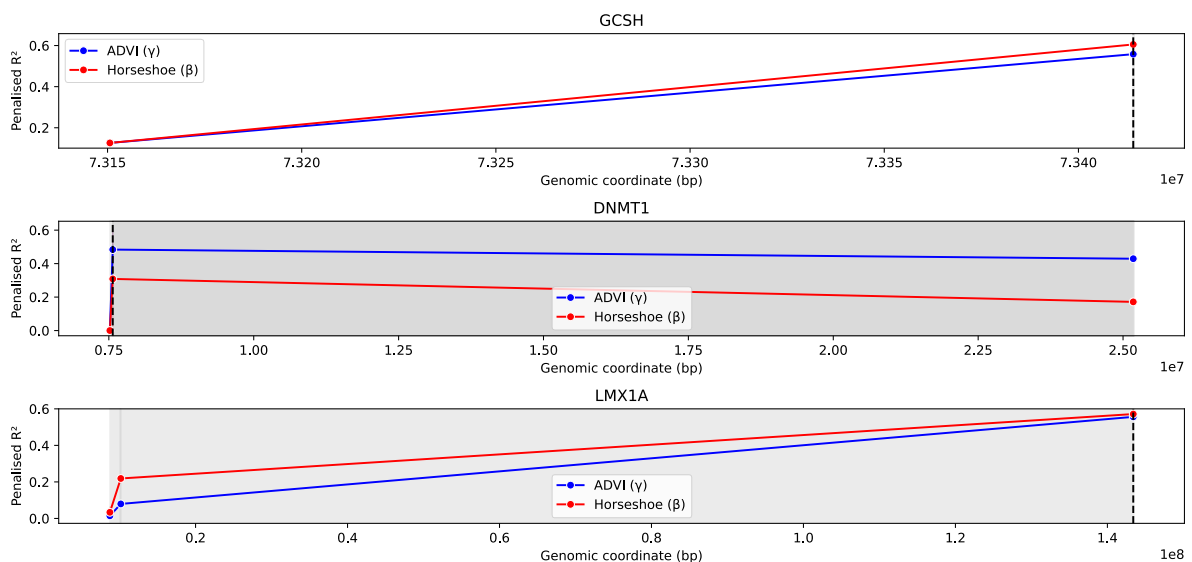
907

908
$$H_{k,i} = \tilde{v}_{\pi(i)}^{(k)}$$

909 **1.12 Data and code availability:**

910 All data and code is available in the study's repository: ...

911 **Supplementary Figure 1: The sliding window self-tuning mechanism was illustrated to**
 912 **the three genes determined by four-CpG windows.** The three genes were mapped for the
 913 genomic coordinates of *Z. californianus*. The coloured dots represent the individual CpGs, the
 914 black dashed line represents the termination of the biological signal (difference in directionality
 915 of the effect size). The grey shading represents the regions for which there is a consistent
 916 change in direction for both the ADVI and Horseshoe models. The genomic coordinates/ their
 917 spans are illustrated on the x axis. The genomic ranges for the genes are: 25Kb (GCSH), 31.7Kb
 918 (DNMT1) and 5.46Kb (LMX1A).



919

920 **References:**

- 921 1. Verplaetse N, Passemiers A, Arany A, Moreau Y, Raimondi D. Large sample size
 922 and nonlinear sparse models outline epistatic effects in inflammatory bowel disease.
 923 Genome Biol. 2023;24(1):224.
- 924 2. Kirpich A, Ainsworth EA, Wedow JM, Newman JRB, Michailidis G, McIntyre LM.
 925 Variable selection in omics data: A practical evaluation of small sample sizes. PLoS One.
 926 2018;13(6):e0197910.

- 927 3. Flesch EP, Rotella JJ, Thomson JM, Graves TA, Garrott RA. Evaluating sample size
928 to estimate genetic management metrics in the genomics era. *Molecular Ecology*
929 *Resources*. 2018;18(5):1077-91.
- 930 4. Vicente-Santos A, Sandoval-Herrera N, Czirják GÁ, Neely BA, Becker DJ.
931 *Proteomics Approaches to Ecoimmunology: New Insights into Wildlife Immunity and*
932 *Disease. Integrative and Comparative Biology*. 2025;65(6):1744-54.
- 933 5. Bogan SN, Yi SV. Potential Role of DNA Methylation as a Driver of Plastic
934 Responses to the Environment Across Cells, Organisms, and Populations. *Genome*
935 *Biology and Evolution*. 2024;16(2):evae022.
- 936 6. Robeck TR, Haghani A, Fei Z, Lindemann DM, Russell J, Herrick KES, et al. Multi-
937 tissue DNA methylation aging clocks for sea lions, walruses and seals. *Commun Biol*.
938 2023;6(1):359.
- 939 7. Photopoulou T, Durbach I, Pirota E, Barratclough A, Schwacke LH, Takeshita R, et
940 al. Methods for analysing wildlife DNA methylation data. *Conserv Physiol*.
941 2026;14(1):coaf091.
- 942 8. Uffelmann E, Huang QQ, Munung NS, de Vries J, Okada Y, Martin AR, et al.
943 *Genome-wide association studies. Nature Reviews Methods Primers*. 2021;1(1):59.
- 944 9. Pei Y, Li G. Unraveling the interplay of DNA methylation and chromosome
945 organization. *Biochem Soc Trans*. 2025;53(6):1433-44.
- 946 10. Weise MJ, Costa DP. Total body oxygen stores and physiological diving capacity of
947 California sea lions as a function of sex and age. *J Exp Biol*. 2007;210(Pt 2):278-89.
- 948 11. Fowler SL, Costa DP, Arnould JP, Gales NJ, Kuhn CE. Ontogeny of diving behaviour
949 in the Australian sea lion: trials of adolescence in a late bloomer. *J Anim Ecol*.
950 2006;75(2):358-67.

- 951 12. McInnes L, Healy J, Astels S. hdbscan: Hierarchical density based clustering. The
952 Journal of Open Source Software. 2017;2:205.
- 953 13. Fisher RA. THE USE OF MULTIPLE MEASUREMENTS IN TAXONOMIC PROBLEMS.
954 Annals of Eugenics. 1936;7(2):179-88.
- 955 14. Kucukelbir A, Tran D, Ranganath R, Gelman A, Blei D. Automatic Differentiation
956 Variational Inference. 2016;18.
- 957 15. Carvalho CM, Polson NG, Scott JG. The horseshoe estimator for sparse signals.
958 Biometrika. 2010;97(2):465-80.
- 959 16. McLean CY, Bristor D, Hiller M, Clarke SL, Schaar BT, Lowe CB, et al. GREAT
960 improves functional interpretation of cis-regulatory regions. Nat Biotechnol.
961 2010;28(5):495-501.
- 962 17. Kruskal JB. On the shortest spanning subtree of a graph and the traveling
963 salesman problem. Proceedings of the American Mathematical society. 1956;7(1):48-50.
- 964 18. Haghverdi L, Buttner M, Wolf FA, Buettner F, Theis FJ. Diffusion pseudotime
965 robustly reconstructs lineage branching. Nat Methods. 2016;13(10):845-8.
- 966 19. Thomas K, Harvey JT, Goldstein T, Barakos J, Gulland F. Movement, dive behavior,
967 and survival of California sea lions (*Zalophus californianus*) posttreatment for domoic
968 acid toxicosis. Marine Mammal Science. 2010;26(1):36-52.
- 969 20. Crofts SJC, Latorre-Crespo E, Chandra T. DNA methylation rates scale with
970 maximum lifespan across mammals. Nat Aging. 2024;4(1):27-32.
- 971 21. Robin ED, Ensinck J, Hance AJ, Newman A, Lewiston N, Cornell L, et al.
972 Glucoregulation and simulated diving in the harbor seal *Phoca vitulina*. Am J Physiol.
973 1981;241(5):R293-300.

- 974 22. Hance AJ, Robin ED, Halter JB, Lewiston N, Robin DA, Cornell L, et al. Hormonal
975 changes and enforced diving in the harbor seal *Phoca vitulina*. II. Plasma
976 catecholamines. *Am J Physiol*. 1982;242(5):R528-32.
- 977 23. Davis RW. Lactate and glucose metabolism in the resting and diving harbor seal
978 (*Phoca vitulina*). *Journal of comparative physiology*. 1983;153(2):275-88.
- 979 24. Arribas-Carreira L, Dallabona C, Swanson MA, Farris J, Ostergaard E, Tsiakas K, et
980 al. Pathogenic variants in GCSH encoding the moonlighting H-protein cause combined
981 nonketotic hyperglycinemia and lipoate deficiency. *Hum Mol Genet*. 2023;32(6):917-33.
- 982 25. Majethia P, Somashekar PH, Hebbar M, Kadavigere R, Praveen BK, Girisha KM, et
983 al. Biallelic start loss variant, c.1A > G in GCSH is associated with variant nonketotic
984 hyperglycinemia. *Clin Genet*. 2021;100(2):201-5.
- 985 26. Marinkovich M, Wack RF, Field CL, Whoriskey ST, Kass PH, Gjeltrema J. Evaluation
986 of Serial Blood Lactate and the Use of a Point-of-Care Lactate Meter in Live-Stranded
987 Pinnipeds. *J Zoo Wildl Med*. 2019;50(1):137-46.
- 988 27. Huang T, Zhou X, Mao X, Yu C, Zhang Z, Yang J, et al. Lactate-fueled oxidative
989 metabolism drives DNA methyltransferase 1-mediated transcriptional co-activator with
990 PDZ binding domain protein activation. *Cancer Sci*. 2020;111(1):186-99.
- 991 28. Xu R, Sun Y, Chen Z, Yao Y, Ma G. Hypoxic Preconditioning Inhibits Hypoxia-
992 induced Apoptosis of Cardiac Progenitor Cells via the PI3K/Akt-DNMT1-p53 Pathway. *Sci*
993 *Rep*. 2016;6:30922.
- 994 29. Paz MF, Wei S, Cigudosa JC, Rodriguez-Perales S, Peinado MA, Huang TH, et al.
995 Genetic unmasking of epigenetically silenced tumor suppressor genes in colon cancer
996 cells deficient in DNA methyltransferases. *Hum Mol Genet*. 2003;12(17):2209-19.

- 997 30. Doucet-Beaupre H, Gilbert C, Profes MS, Chabrat A, Pacelli C, Giguere N, et al.
998 Lmx1a and Lmx1b regulate mitochondrial functions and survival of adult midbrain
999 dopaminergic neurons. *Proc Natl Acad Sci U S A*. 2016;113(30):E4387-96.
- 1000 31. Xie J, Cai TT, Li H. Sample size and power analysis for sparse signal recovery in
1001 genome-wide association studies. *Biometrika*. 2011;98(2):273-90.
- 1002 32. Forsman A. Effects of genotypic and phenotypic variation on establishment are
1003 important for conservation, invasion, and infection biology. *Proc Natl Acad Sci U S A*.
1004 2014;111(1):302-7.
- 1005 33. Goldberg DC, Fu H, Atkins D, Moyer E, Lee CN, Deng Y, et al. KnowYourCG:
1006 Facilitating base-level sparse methylome interpretation. *Sci Adv*. 2025;11(43):eadw3027.
- 1007 34. Love Stowell SM, Bentley EG, Gagne RB, Gustafson KD, Rutledge LY, Ernest HB.
1008 Optimal DNA extractions from blood on preservation paper limits conservation genomic
1009 but not conservation genetic applications. *Journal for Nature Conservation*. 2018;46:89-
1010 96.
- 1011 35. Esipova P, Suvorova I, Yachmen V, Pushchin I. Blood Morphology and Hematology
1012 of Adult Baikal Seals (*Pusa sibirica* Gmelin, 1788) Under Professional Care. *Animals*
1013 (Basel). 2025;15(2).
- 1014 36. Mizoguchi BA, Valenzuela N. A Cautionary Tale of Sexing by Methylation: Hybrid
1015 Bisulfite-Conversion Sequencing of Immunoprecipitated Methylated DNA in *Chrysemys*
1016 *picta* Turtles with Temperature-Dependent Sex Determination Reveals Contrasting
1017 Patterns of Somatic and Gonadal Methylation, but No Unobtrusive Sex Diagnostic.
1018 *Animals (Basel)*. 2022;13(1).
- 1019