

000 001 002 003 004 005 006 007 008 009 010 011 012 013 014 015 016 017 018 019 020 021 022 023 024 025 026 027 028 029 030 031 032 033 034 035 036 037 038 039 040 041 042 043 044 045 046 047 048 049 050 051 052 053 GEOBIS: BUDGET-OPTIMAL BLOCK IMPORTANCE SAMPLING FOR STOCHASTIC RIEMANNIAN OPTIMI- ZATION

006
007 **Anonymous authors**
008 Paper under double-blind review

010 011 ABSTRACT

013 This paper studies budgeted block subsampling for stochastic Riemannian optimiza-
014 tion. Starting from the Horvitz–Thompson estimator, we derive an independent
015 Bernoulli design with a water-filling probability rule that minimizes the second
016 moment under a fixed expected number of active blocks. The resulting estimator,
017 GeoBIS, is unbiased and achieves the canonical inverse-in-budget behavior of its
018 second moment. We also analyze exact- K negatively dependent designs, includ-
019 ing projection determinantal point processes and sampling without replacement
020 with unequal probabilities. Under a mild alignment condition on block directions,
021 exact- K strictly reduces the cross term in the variance. A simple wall-clock model
022 provides a closed-form rule for selecting the active-block budget and clarifies
023 when exact- K is worthwhile. Experiments on orthogonality-constrained sequence
024 models and thin-Stiefel adapters follow the predicted trends and validate GeoBIS
025 as a practical default.

026 027 1 INTRODUCTION

028
029 Optimization with manifold constraints appears in many modern machine learning systems. Examples
030 include orthogonal and Stiefel constraints in sequence models and adapters, subspace learning and
031 matrix factorizations, and modules on product or hyperbolic manifolds. In these settings, the
032 Riemannian gradient lives in a tangent space that admits natural block decompositions: coordinate
033 groups, skew-rotation atoms such as Givens or Householder generators, or per-factor blocks in
034 product manifolds. Executing full tangent updates can be expensive, and a common engineering
035 strategy is to activate only a subset of blocks at each step and to reweight the contribution of sampled
036 blocks using Horvitz–Thompson scaling so that the gradient estimator remains unbiased.

037 This paper addresses the design of such block subsampling policies when a fixed compute budget
038 limits the expected number of active blocks per step. The first question is whether one can choose
039 marginal probabilities so that the resulting estimator minimizes its second moment for a given budget.
040 The second question is whether using a fixed-size subset with negative dependence can further reduce
041 redundancy among sampled blocks when their directions are aligned. The third question is how to
042 choose the budget itself in a way that optimizes wall-clock time to a target accuracy, while keeping
043 the method practical.

044 The contributions are organized as follows. First, we present GeoBIS, an independent Bernoulli
045 design that minimizes the Horvitz–Thompson second moment under a fixed expected number of
046 active blocks through a simple water-filling rule. The formula admits an exact capped variant when
047 some probabilities hit one. Second, we give a design-agnostic lower bound showing that the inverse-
048 in-budget dependence of the dominant term in the second moment is unavoidable for any unbiased
049 design; this identifies a statistical frontier and explains why water-filling is instance-optimal up to
050 constants. Third, we study exact- K negatively dependent designs, including projection determinantal
051 point processes and unequal-probability sampling without replacement, and we state a precise
052 alignment condition under which these designs strictly reduce the cross term in the variance compared
053 to independent Bernoulli sampling with matched marginals. Fourth, we provide a wall-clock analysis
054 with a simple cost model that yields a closed-form rule for setting the active-block budget and an
055 online schedule based on exponential moving averages. Finally, we discuss orthogonal and Stiefel

054 instantiations and present experiments whose focus is on validating the statistical and compute trends
 055 rather than on competitive accuracy.
 056

057 The paper is written to separate essential ideas from optional extensions. Each section begins with
 058 context and design choices before introducing the formal statements.
 059

060 2 RELATED WORK 061

062 Riemannian optimization is by now a standard tool for constrained machine learning. Foundational
 063 references include treatments of retractions, vector transports, and convergence guarantees for first-
 064 order methods and trust-region methods. Monographs such as Absil et al. (2008) and Boumal (2023)
 065 cover the relevant geometric preliminaries and provide a unifying view of manifold-constrained algo-
 066 rithms. Early stochastic Riemannian methods and their convergence properties appear in Bonnabel
 067 (2013), while analyses for geodesically convex problems and smoothness under retractions are studied
 068 in Zhang & Sra (2016). Variance reduction in Riemannian settings has been explored in Sato et al.
 069 (2019) and quasi-Newton adaptations in Kasai et al. (2018).
 070

071 Orthogonality constraints arise in several model classes. For feasible updates with orthogonality
 072 constraints, Cayley-transform based methods and related retractions are discussed in Wen & Yin
 073 (2013). Their use within deep learning layers and flows appears in Lezcano-Casado (2019) and
 074 Trockman & Kolter (2021). These works focus on preserving constraints efficiently and motivate
 075 block choices aligned with the geometry, such as skew-symmetric generators for the orthogonal group
 076 and projected gradients on the Stiefel manifold.
 077

078 In Euclidean settings, randomized coordinate descent and block coordinate methods have a long
 079 history. Analyses that allow arbitrary sampling and show the role of importance sampling include
 080 Nesterov (2012); Richtárik & Takáč (2016); Qu et al. (2015). Without-replacement sampling and
 081 random reshuffling often improve practical convergence behavior compared to with-replacement
 082 sampling, as discussed in Shamir (2016); Gürbüzbalaban et al. (2021); HaoChen & Sra (2019). The
 083 present paper brings budget-explicit sampling design and compute-aware analysis to Riemannian
 084 block updates, where geometric alignment enables sharper variance formulas and practical selection
 085 rules.
 086

087 Negative dependence and diversity-seeking fixed-size designs can reduce redundancy in subsampled
 088 sets. Determinantal point processes provide a tractable family of negatively associated distributions
 089 over subsets with closed-form marginals and pairwise inclusion probabilities; see Kulesza & Taskar
 090 (2012); Hough et al. (2006). Projection determinantal point processes, built from an orthogonal
 091 projector, are particularly convenient because they maintain a fixed size and admit simple sampling
 092 routines once the top eigenspace is available. The use of such designs for randomized numerical
 093 linear algebra and column subset selection is surveyed in Dereziński & Mahoney (2020), with related
 094 ideas in volume sampling Deshpande et al. (2006) and leverage-score based methods Cohen et al.
 095 (2017). In survey sampling, classical schemes for unequal-probability sampling without replacement
 096 include Sampford (1967) and pivotal splitting methods such as Deville & Tillé (1998), which avoid
 097 eigen-decompositions while matching prescribed marginals. We use these designs as optional exact-
 098 K extensions of the independent sampler and clarify the condition under which they strictly reduce
 099 the variance.
 100

101 Finally, recent work on Riemannian coordinate or block methods with iteration complexity guarantees
 102 includes Han et al. (2024). Our approach differs by focusing on the design of unbiased estimators
 103 with minimal second moment under an explicit compute budget, by giving closed-form rules for
 104 marginal probabilities and for the active-block budget itself, and by quantifying when and why
 105 negative dependence helps in practice.
 106

107 3 SETUP AND NOTATION 108

109 This section sets the notation and states the estimator family under study. The goal is to mini-
 110 mize a smooth objective on a Riemannian manifold using unbiased stochastic gradients formed by
 111 subsampling blocks of the tangent space.
 112

108 Let (\mathcal{M}, g) be a finite-dimensional Riemannian manifold. We minimize a smooth function $F : \mathcal{M} \rightarrow$
 109 \mathbb{R} using retraction-based updates of the form
 110

$$111 \quad X_{t+1} = \text{Retr}_{X_t} \left(-\eta_t \hat{\xi}_t \right), \quad \hat{\xi}_t \in T_{X_t} \mathcal{M}, \quad (1)$$

113 where $\eta_t > 0$ is a stepsize, $\hat{\xi}_t$ is an unbiased estimator of the Riemannian gradient $\xi_t = \text{grad}F(X_t)$,
 114 and Retr is a first-order accurate retraction.

115 At a point X , assume a block decomposition of the tangent space
 116

$$117 \quad T_X \mathcal{M} = \bigoplus_{b=1}^B V_b, \quad \xi = \sum_{b=1}^B \xi_{(b)}, \quad \xi_{(b)} \in V_b. \quad (2)$$

120 Define block magnitudes $v_b = \|\xi_{(b)}\|_{g_X}$ and unit block directions $u_b = \xi_{(b)}/v_b$ for $v_b > 0$. The
 121 directional cosines between blocks are

$$122 \quad \rho_{bc} = \langle u_b, u_c \rangle_{g_X} \in [-1, 1], \quad G = [\rho_{bc}]. \quad (3)$$

123 We allow mild deviations from exact block orthogonality and quantify them by a coherence parameter
 124 $\mu = \max_{b \neq c} |\rho_{bc}|$.

126 The estimator family is based on Horvitz–Thompson scaling. Let $S \subset [B]$ be a random subset
 127 of blocks with first-order inclusion probabilities $\pi_b = \Pr(b \in S)$ and second-order inclusion
 128 probabilities $\pi_{bc} = \Pr(\{b, c\} \subset S)$. The estimator is

$$129 \quad \hat{\xi} = \sum_{b \in S} \frac{\xi_{(b)}}{\pi_b}. \quad (4)$$

132 It is unbiased because $\mathbb{E}[\hat{\xi}] = \sum_b \Pr(b \in S) \xi_{(b)}/\pi_b = \sum_b \xi_{(b)} = \xi$. The second moment and
 133 variance decompose as follows.

134 Lemma 1 (Variance decomposition). For any unbiased Horvitz–Thompson estimator supported on a
 135 block subset S ,

$$137 \quad \text{Var}(\hat{\xi}) = \sum_b \left(\frac{1}{\pi_b} - 1 \right) v_b^2 + 2 \sum_{b < c} \left(\frac{\pi_{bc}}{\pi_b \pi_c} - 1 \right) v_b v_c \rho_{bc}. \quad (5)$$

139 The first term depends only on the marginals and is the dominant quantity to optimize under a
 140 budget. The second term depends on dependence across block indicators; it is zero under independent
 141 Bernoulli sampling and becomes negative under negatively associated exact- K designs when the
 142 block directions are aligned.

144 4 GEOBIS: BUDGET-OPTIMAL BERNOUlli VIA WATER-FILLING

146 This section presents the independent design. Each block b is included independently with probability
 147 $p_b \in (0, 1]$ and is scaled by I_b/p_b . The expected number of active blocks is the budget $K = \sum_b p_b$.
 148 When blocks are orthogonal in the metric, the expected second moment is $\sum_b v_b^2/p_b$. We now
 149 minimize this quantity under the linear budget constraint and recover a water-filling rule.

151 Proposition 1 (Water-filling). Consider the strictly convex problem $\min_{p \in (0, 1]^B} \sum_b v_b^2/p_b$ subject to
 152 $\sum_b p_b = K$. The unique solution is

$$153 \quad p_b^* = \min\{1, \lambda v_b\}, \quad \text{where } \lambda > 0 \text{ is chosen so that } \sum_b p_b^* = K. \quad (6)$$

156 In the no-cap regime $p_b < 1$ for all b , this yields $p_b^* = K v_b / \sum_j v_j$ and the expected second moment
 157 equals $(\sum_b v_b)^2/K$.

159 The proof uses KKT conditions with inequality constraints $p_b \leq 1$ and the convexity of $x \mapsto 1/x$.
 160 An efficient implementation sorts the v_b once per step to determine the capped set and computes
 161 the threshold λ by a one-pass sweep. When caps are active, the second moment decomposes into a
 162 capped part plus an uncapped part divided by the residual budget.

Practical safety. Horvitz–Thompson scaling uses $1/p_b$ and can be numerically unstable if a proxy misclassifies a very small v_b . We recommend a probability floor p_{\min} and use $p_b = \min\{1, \max\{p_{\min}, \lambda v_b\}\}$. We also recommend smoothing scores with exponentials moving averages to reduce noise.

Mild non-orthogonality. When blocks are nearly orthogonal, the Bernoulli formulas remain accurate. The following bound expresses how much the variance can deviate from the independent-marginal term.

Proposition 2 (Coherence-based robustness). Let $\mu = \max_{b \neq c} |\rho_{bc}|$. For any exact- K design with first- and second-order inclusions (π_b, π_{bc}) ,

$$\left| \text{Var}(\hat{\xi}) - \sum_b \left(\frac{1}{\pi_b} - 1 \right) v_b^2 \right| \leq 2\mu \sum_{b < c} \left| \frac{\pi_{bc}}{\pi_b \pi_c} - 1 \right| v_b v_c. \quad (7)$$

In particular, under Bernoulli sampling the right-hand side is zero, and under negatively associated exact- K designs it is controlled by the size of off-diagonal entries.

Proxy scores and regret. In practice one may water-fill on proxy scores s_b that approximate v_b within a multiplicative error $(1 \pm \epsilon)$, and occasionally misidentify the top capped set. In the no-cap regime, the multiplicative inflation of the second moment is at most $(1 + \epsilon)/(1 - \epsilon)$. With caps and a misranking probability δ of the top set, the excess second moment is bounded by a term proportional to $\delta(1 + \epsilon)/(1 - \epsilon)$ times the uncapped contribution.

5 EXACT- K NEGATIVE-DEPENDENCE EXTENSIONS

This section discusses fixed-size designs that discourage redundant co-selection of aligned blocks. We present projection determinantal point processes and unequal-probability sampling without replacement, state the variance formula under exact- K , and give a condition under which exact- K strictly improves over Bernoulli at matched marginals.

Projection determinantal point processes. Let $W = \text{diag}(\sqrt{v_b})$ and $G = [\rho_{bc}]$. Form the geometry-weighted Gram matrix $M = W G W$. Let V_K be the top- K eigenspace of M and $K^\sharp = V_K V_K^\top$ the projection kernel. The projection k -DPP with kernel K^\sharp satisfies $\pi_b = K_{bb}^\sharp$ and $\pi_{bc} = \pi_b \pi_c - (K_{bc}^\sharp)^2$. By construction, the subset size is K , and the block indicators are negatively associated.

Unequal-probability sampling without replacement. When one desires exact- K without eigenspace computations, pivotal Poisson and related Sampford-type schemes provide fixed-size sampling with prescribed marginals π_b . These designs maintain Horvitz–Thompson unbiasedness and avoid spectral preprocessing.

Variance comparison at matched marginals. The variance of the Horvitz–Thompson estimator under exact- K reads

$$\text{Var}(\hat{\xi}) = \sum_b \left(\frac{1}{\pi_b} - 1 \right) v_b^2 + 2 \sum_{b < c} \left(\frac{\pi_{bc}}{\pi_b \pi_c} - 1 \right) v_b v_c \rho_{bc}. \quad (8)$$

For the projection k -DPP, the pairwise term becomes $-2 \sum_{b < c} (K_{bc}^\sharp)^2 v_b v_c \rho_{bc} / (\pi_b \pi_c)$. A direct comparison to Bernoulli with the same marginals therefore requires a condition that this signed sum is nonnegative. We phrase this as an aggregate alignment condition.

Assumption 1 (Aggregate alignment). Let $\pi_b = K_{bb}^\sharp$ and define weights $w_{bc} = (K_{bc}^\sharp)^2 v_b v_c / (\pi_b \pi_c)$. Assume that $\sum_{b < c} w_{bc} \rho_{bc} \geq 0$. This holds, for example, if all $\rho_{bc} \geq 0$ on the support of K^\sharp .

Theorem 1 (Exact- K variance reduction under alignment). Let $\hat{\xi}_{\text{DPP}}$ be the Horvitz–Thompson estimator under the projection k -DPP with kernel K^\sharp and marginals $\pi_b = K_{bb}^\sharp$. Let $\hat{\xi}_{\text{Bern}}$ be the Bernoulli estimator with the same marginals. Under the aggregate alignment assumption,

$$\text{Var}(\hat{\xi}_{\text{DPP}}) \leq \text{Var}(\hat{\xi}_{\text{Bern}}), \quad (9)$$

216 with strict inequality if there exists $b \neq c$ such that $\rho_{bc} > 0$ and $K_{bc}^\sharp \neq 0$.
 217

218 This result clarifies when negative dependence is beneficial. On manifolds and block choices where
 219 cross-block cosines are small, such as disjoint Givens generators on the orthogonal group, the
 220 improvement is minor in practice and the extra selection overhead can outweigh the benefit. On
 221 models with coherent blocks, exact- K can reduce the variance at matched marginals.
 222

223 6 LOWER BOUND AND INSTANCE OPTIMALITY

225 We now show that the inverse-in-budget dependence of the dominant term in the second moment
 226 is unavoidable for any unbiased estimator that respects the budget on expected active blocks. This
 227 identifies a statistical frontier and explains why the independent water-filling design is instance-
 228 optimal up to constants.
 229

230 Theorem 2 (Design-agnostic lower bound). Let $\hat{\xi}$ be any unbiased Horvitz–Thompson estimator with
 231 marginals summing to K . Then

$$232 \mathbb{E}\|\hat{\xi}\|^2 \geq \|\xi\|^2 + \frac{(\sum_b v_b)^2}{K} - (\sum_b v_b)^2. \quad (10)$$

235 Equivalently, the variance satisfies $\text{Var}(\hat{\xi}) \geq (\sum_b v_b)^2/K - \sum_b v_b^2$.
 236

237 The proof bounds the cross term in the variance decomposition below and uses Cauchy–Schwarz to
 238 minimize $\sum_b v_b^2/\pi_b$ subject to $\sum_b \pi_b = K$. The independent design with water-filling achieves the
 239 lower bound in the no-cap regime. Exact- K designs can only subtract a nonnegative cross term when
 240 the alignment condition holds.
 241

242 7 NONCONVEX PROGRESS AND WALL-CLOCK SCHEDULING

244 We briefly recall a standard retraction-smoothness descent lemma and combine it with the second-
 245 moment formulas to motivate a simple rule for choosing the active-block budget.
 246

Assumption 2 (Retraction smoothness). There exists $L > 0$ such that for all X and small $\zeta \in T_X \mathcal{M}$,

$$248 F(\text{Retr}_X(\zeta)) \leq F(X) + \langle \text{grad}F(X), \zeta \rangle_{g_X} + \frac{L}{2} \|\zeta\|_{g_X}^2. \quad (11)$$

250 Under unbiasedness and bounded second moment $\mathbb{E}\|\hat{\xi}_t\|^2 \leq G_t$, constant stepsize $\eta \leq 1/L$ yields
 251

$$252 \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}\|\text{grad}F(X_t)\|^2 \leq \frac{2(F(X_0) - F^{\inf})}{\eta T} + L \eta \bar{G}, \quad \bar{G} = \frac{1}{T} \sum_t G_t. \quad (12)$$

255 For independent water-filling without caps, a coarse but useful model writes $G_t \approx a_t + b_t/K$
 256 with $a_t = \sum_b v_{t,b}^2$ and $b_t = (\sum_b v_{t,b})^2$. If the per-step cost is $C_0 + C_1 K$ and we approximate a_t
 257 and b_t by slowly varying averages \bar{a} and \bar{b} , then the time to a target stationarity is proportional to
 258 $(C_0 + C_1 K)(\bar{a} + \bar{b}/K)$. The minimizer is
 259

$$260 K^* = \sqrt{\frac{C_0 \bar{b}}{C_1 \bar{a}}}. \quad (13)$$

262 This rule is myopic and assumes local stationarity of \bar{a} and \bar{b} . In practice, one tracks exponential
 263 moving averages and clips K^* to $[1, B]$. Under exact- K designs that satisfy the alignment condition,
 264 the second moment is further reduced by the negative cross term, while the selection overhead
 265 increases the fixed per-step cost. The same balancing logic applies.
 266

267 8 ORTHOGONAL AND STIEFEL INSTANTIATIONS

268 This section describes the geometry-aware block choices and local formulas used in practice.
 269

270 **Orthogonal group.** The orthogonal group $O(d)$ consists of matrices $X \in \mathbb{R}^{d \times d}$ with $X^\top X = I$.
 271 The tangent space at X is $\{Z : X^\top Z + Z^\top X = 0\}$. For an Euclidean gradient $G = \nabla_X F$, the
 272 Riemannian gradient under the canonical metric is

$$275 \quad \text{grad}F(X) = G - X \text{ sym}(X^\top G) = \Omega X, \quad \Omega = \frac{1}{2} (GX^\top - XG^\top). \quad (14)$$

278 A natural block family is given by skew-symmetric generators $E_{ij} - E_{ji}$. The block magnitudes are
 279 $v_{ij} = \sqrt{2} |\Omega_{ij}|$, and the water-filling rule uses these local scores. Retractions include Cayley and
 280 exponential maps on skew subspaces; disjoint pairs update in parallel.

284 **Stiefel manifold.** The thin Stiefel manifold $\text{St}(d, p)$ consists of matrices $X \in \mathbb{R}^{d \times p}$ with $X^\top X =$
 285 I_p . The Riemannian gradient is the projected Euclidean gradient $G - X \text{ sym}(X^\top G)$. Natural
 286 blocks are row- or column-groups aligned with the metric. Retractions include QR- and polar-based
 287 retractions and Cayley-type updates adapted to the constraint.

290 9 EXPERIMENTS

293 The experiments are designed to validate the statistical and compute trends predicted by the analysis
 294 rather than to optimize task accuracy. We consider three sequence datasets with orthogonality-
 295 constrained recurrent models and a thin Stiefel adapter inserted in a small convolutional network.
 296 The methods compared are independent water-filling and uniform sampling at matched budget and a
 297 full-gradient baseline. Negative-dependence exact- K designs are disabled in these regimes because
 298 selection overhead dominates.

299 Tables summarize wall-clock time, validation accuracy, and the mean second moment of the
 300 Horvitz–Thompson estimator. The inverse-in-budget trend for the second moment under water-
 301 filling is consistent across tasks. Uniform sampling yields larger second moments under the same
 302 budget, as expected from the lower bound.

305 **Orthogonality-constrained sequence models.** Copy, Adding, and psMNIST with a Cayley-
 306 retracted orthogonal RNN. Methods: full gradient, uniform, and GeoBIS with budgets in a small grid.
 307 Metrics: mean second moment, final validation accuracy, and time to reach a fixed fraction of the
 308 best accuracy.

312 Table 1: Copy task with an orthogonality-constrained RNN (Cayley retraction). Budget K is the
 313 expected number of active blocks for sampling-based methods.

315 Method	Budget	Time (s)	Val. Acc. (final)	Mean $\mathbb{E} \ \hat{\xi}\ _g^2$
316 Full gradient	–	319.3	0.1775	1.42e-01
317 GeoBIS (Bernoulli)	4	325.8	0.1775	6.34e+02
318 GeoBIS (Bernoulli)	8	322.7	0.1775	3.21e+02
319 GeoBIS (Bernoulli)	16	324.0	0.1775	1.52e+02
320 Uniform baseline ref	–	321.1	0.1775	1.42e-01
321 Uniform	4	321.7	0.1775	1.24e+03
322 Uniform	8	317.8	0.1775	6.05e+02
323 Uniform	16	320.9	0.1775	2.74e+02

324

325

326 Table 2: Adding task with an orthogonality-constrained RNN (Cayley retraction).

327

328

329

330

331

332

333

334

335

336

337

338

Method	Budget	Time (s)	Val. Acc. (final)	Mean $\mathbb{E}\ \hat{\xi}\ _g^2$
Full gradient	—	283.4	0.5320	2.43e-01
GeoBIS (Bernoulli)	4	309.7	0.5320	1.04e+03
GeoBIS (Bernoulli)	8	305.6	0.5320	5.57e+02
GeoBIS (Bernoulli)	16	310.9	0.5320	2.73e+02
Uniform baseline ref	—	304.1	0.5320	2.43e-01
Uniform	4	301.9	0.5320	2.30e+03
Uniform	8	304.8	0.5320	1.19e+03
Uniform	16	305.0	0.5320	4.95e+02

336

337

338

339 Table 3: psMNIST with an orthogonality-constrained RNN (Cayley retraction). \dagger : run completed
only one epoch, so time and means are not directly comparable to 5-epoch runs.

340

341

342

343

344

345

346

347

348

349

350

351

Method	Budget	Time (s)	Val. Acc. (final)	Mean $\mathbb{E}\ \hat{\xi}\ _g^2$
Full gradient	—	1784.4	0.1256	2.76e+00
GeoBIS (Bernoulli)	4	1239.5	0.0996	1.09e+04
GeoBIS (Bernoulli)	8	1266.7	0.1119	6.12e+03
GeoBIS (Bernoulli) \dagger	16	252.2	0.0980	2.79e+03
Uniform baseline ref	—	1636.7	0.1256	2.76e+00
Uniform	4	1238.8	0.0833	1.96e+04
Uniform	8	1224.3	0.1051	9.04e+03
Uniform	16	1222.9	0.1126	4.25e+03

352

Stiefel adapter. A thin Stiefel adapter after a small convolutional backbone trained on CIFAR-10. Methods and metrics mirror the sequence setting. The observed second-moment reductions for GeoBIS at fixed budget align with the theory, and accuracy differences are small, as expected when backprop dominates compute.

356

357

358 Table 4: Thin Stiefel adapter on CIFAR-10. GeoBIS reduces the mean HT second moment at fixed
budget; accuracy and wall-clock are comparable to Uniform.

359

360

361

362

363

364

365

366

367

368

369

370 Table 5: Thin Stiefel adapter on CIFAR-10: time to reach $0.95 \times$ the best validation accuracy across
all runs.

371

372

373

374

375

376

377

Method	Budget	Time to target (s)
GeoBIS (Bernoulli)	4	228.1
Uniform	4	188.1
GeoBIS (Bernoulli)	8	200.2
Uniform	8	215.7
GeoBIS (Bernoulli)	16	218.6
Uniform	16	213.6

378 **Summary.** Across both domains, water-filling improves the mean second moment at fixed budget
 379 relative to uniform sampling, with trends matching the inverse-in-budget law. Exact- K designs are
 380 left as an opt-in for heavier settings with coherent blocks and lower relative overhead.
 381

382 10 REPRODUCIBILITY STATEMENT 383

384 We will release code (after the review process is complete) that exactly reproduces every number in
 385 the paper’s tables from a clean environment. The package includes: (i) reference implementations
 386 of water-filling, independent sampling with a probability floor, and orthogonal/Stiefel retractions;
 387 (ii) task-specific training and evaluation scripts for the Copy, Adding, psMNIST, and CIFAR-10
 388 thin-Stiefel adapter setups; (iii) configuration files that correspond one-to-one with each table row
 389 (budget K , seeds, optimizer and schedule, batch sizes, retraction type); (iv) and an environment
 390 specification.

391 As for LLM usage, LLMs were used to provide help with paper writing and polishing (especially
 392 around creating table structures, relevant bibtex references, and template fitting).

394 11 LIMITATIONS AND SOCIETAL IMPACT 395

396 The main limitation is that exact- K designs require either eigenspace computations or pivotal
 397 sampling; when the fixed overhead is large and block alignment is weak, the marginal benefit can
 398 be negative. The wall-clock rule is myopic and assumes locally stationary statistics for the block
 399 magnitudes; this is a practical compromise and works well with moving averages but does not
 400 constitute a global optimality result. This work is methodological; we foresee standard concerns
 401 around compute and energy use during training but no special risks beyond those encountered in
 402 typical optimization research.

404 12 CONCLUSION 405

406 GeoBIS provides a simple and budget-optimal independent design for block subsampling in stochastic
 407 Riemannian optimization, with closed-form probabilities and second-moment expressions, including
 408 capped regimes. Exact- K negative-dependence extensions can further reduce variance under an
 409 explicit alignment condition, and a simple cost model yields a closed-form rule for the active-block
 410 budget. The method aligns with the geometry of common manifolds, drops cleanly into existing code,
 411 and is supported by experiments that validate its statistical and compute advantages.

412
 413
 414
 415
 416
 417
 418
 419
 420
 421
 422
 423
 424
 425
 426
 427
 428
 429
 430
 431

432 REFERENCES
433

434 P.-A. Absil, R. Mahony, and R. Sepulchre. Optimization Algorithms on Matrix Manifolds. Princeton
435 University Press, 2008.

436 S. Bonnabel. Stochastic gradient descent on riemannian manifolds. *IEEE Transactions on Automatic
437 Control*, 58(9):2217–2229, 2013.

438

439 N. Boumal. An Introduction to Optimization on Smooth Manifolds. Cambridge University Press,
440 2023. URL <https://www.nicolasboumal.net/book/>.

441 M.B. Cohen, C. Musco, and C. Musco. Input sparsity time low-rank approximation via ridge leverage
442 score sampling. In Proceedings of the Twenty-Eighth Annual ACM-SIAM Symposium on Discrete
443 Algorithms, pp. 1758–1777, 2017.

444

445 M. Dereziński and M.W. Mahoney. Determinantal point processes in randomized numerical linear
446 algebra. arXiv preprint arXiv:2005.03185, 2020.

447 A. Deshpande, L. Rademacher, S. Vempala, and G. Wang. Matrix approximation and projective
448 clustering via volume sampling. *Theory of Computing*, 2(1):225–247, 2006.

449

450 J.-C. Deville and Y. Tillé. Unequal probability sampling without replacement through a splitting
451 method. *Biometrika*, 85(1):89–101, 1998.

452

453 M. Gürbüzbalaban, A. Ozdaglar, and P. Parrilo. Why random reshuffling beats stochastic gradient
454 descent. *Mathematical Programming*, 186(1-2):49–84, 2021.

455

456 J. Han et al. Riemannian coordinate descent methods: iteration complexity and applications. *Mathematical Programming*, 2024. to appear.

457

458 J. HaoChen and S. Sra. Random shuffling beats sgd after finite epochs. In *International Conference
459 on Machine Learning*, 2019.

460

461 J.B. Hough, M. Krishnapur, Y. Peres, and B. Virág. Determinantal processes and independence.
Probability Surveys, 3:206–229, 2006.

462

463 H. Kasai, H. Sato, and B. Mishra. Riemannian stochastic quasi-newton algorithm with variance
464 reduction. In *Proceedings of Machine Learning Research*, volume 84, pp. 1852–1861, 2018.

465

466 A. Kulesza and B. Taskar. Determinantal point processes for machine learning. *Foundations and
Trends in Machine Learning*, 5(2–3):123–286, 2012.

467

468 M. Lezcano-Casado. Trivializations for gradient-based optimization on manifolds. In *Advances in
Neural Information Processing Systems*, 2019.

469

470 Y. Nesterov. Efficiency of coordinate descent methods on huge-scale optimization problems. *SIAM
Journal on Optimization*, 22(2):341–362, 2012.

471

472 Z. Qu, P. Richtárik, and T. Zhang. QUARTZ: Randomized dual coordinate ascent with arbitrary
473 sampling. In *Advances in Neural Information Processing Systems*, 2015.

474

475 P. Richtárik and M. Takáč. Parallel coordinate descent methods for big data optimization. *Mathematical
Programming*, 156(1-2):433–484, 2016.

476

477 M.R. Sampford. On sampling without replacement with unequal probabilities of selection. *Biometrika*,
478 54(3-4):499–513, 1967.

479

480 H. Sato, H. Kasai, and B. Mishra. Riemannian stochastic variance reduced gradient algorithm with
481 retraction and vector transport. *SIAM Journal on Optimization*, 29(4):3130–3162, 2019.

482

483 O. Shamir. Without-replacement sampling for stochastic gradient methods. In *Advances in Neural
Information Processing Systems*, 2016.

484

485 A. Trockman and J.Z. Kolter. Orthogonalizing convolutional layers with the cayley transform.
International Conference on Learning Representations, 2021. arXiv:2104.07167.

486 Z. Wen and W. Yin. A feasible method for optimization with orthogonality constraints. Mathematical
487 Programming, 142(1):397–434, 2013.
488

489 H. Zhang and S. Sra. First-order methods for geodesically convex optimization. In Conference on
490 Learning Theory, pp. 1617–1638, 2016.
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

540 **A ADDITIONAL PRELIMINARIES**
 541

542 This appendix collects definitions and lemmas used in the main text. A retraction $\text{Retr}_X : T_X \mathcal{M} \rightarrow$
 543 \mathcal{M} is a smooth map that satisfies $\text{Retr}_X(0_X) = X$ and whose differential at the origin equals the
 544 identity. Retraction-smoothness states that along retraction curves, the objective is upper-bounded
 545 by a quadratic with constant L . For stochastic analysis, we assume unbiasedness and a bounded
 546 second moment for the estimator. Block decompositions may be only approximately orthogonal; the
 547 coherence parameter μ controls the magnitude of off-diagonal terms in inner products between unit
 548 block directions.

549 **B VARIANCE DECOMPOSITION PROOF**
 550

552 For completeness, we provide a short derivation of the variance identity. Write $\widehat{\xi} - \xi = \sum_b (\frac{I_b}{\pi_b} - 1) \xi_{(b)}$.
 553 Since $\mathbb{E}[\frac{I_b}{\pi_b} - 1] = 0$ and $\mathbb{E}[(\frac{I_b}{\pi_b} - 1)(\frac{I_c}{\pi_c} - 1)] = \frac{\pi_{bc}}{\pi_b \pi_c} - 1$ for $b \neq c$,

$$555 \quad \mathbb{E}\|\widehat{\xi} - \xi\|^2 = \sum_b \mathbb{E}\left(\frac{I_b}{\pi_b} - 1\right)^2 \|\xi_{(b)}\|^2 + 2 \sum_{b < c} \mathbb{E}\left(\frac{I_b}{\pi_b} - 1\right) \left(\frac{I_c}{\pi_c} - 1\right) \langle \xi_{(b)}, \xi_{(c)} \rangle, \quad (15)$$

$$558 \quad \rangle = \sum_b \left(\frac{1}{\pi_b} - 1\right) v_b^2 + 2 \sum_{b < c} \left(\frac{\pi_{bc}}{\pi_b \pi_c} - 1\right) v_b v_c \rho_{bc}. \quad (16)$$

560 This equals the variance because the estimator is unbiased.

562 **C WATER-FILLING WITH CAPS**
 563

564 We solve $\min \sum_b v_b^2 / p_b$ subject to $\sum_b p_b = K$ and $0 < p_b \leq 1$. The Lagrangian is

$$567 \quad \mathcal{L}(p, \mu, \alpha) = \sum_b \frac{v_b^2}{p_b} + \mu \left(\sum_b p_b - K \right) + \sum_b \alpha_b (p_b - 1), \quad (17)$$

569 with $\alpha_b \geq 0$. First-order conditions yield $-\frac{v_b^2}{p_b^2} + \mu + \alpha_b = 0$. If $p_b < 1$, then $\alpha_b = 0$ and
 570 $p_b = \sqrt{v_b^2 / \mu} = \lambda v_b$ with $\lambda = 1 / \sqrt{\mu}$. If $\lambda v_b > 1$, then $p_b = 1$ and the corresponding α_b is
 571 positive. Let $\mathcal{C} = \{b : \lambda v_b \geq 1\}$ and $\mathcal{U} = [B] \setminus \mathcal{C}$ with $|\mathcal{C}| = m$. The budget constraint becomes
 572 $m + \lambda \sum_{b \in \mathcal{U}} v_b = K$, hence $\lambda = (K - m) / S_{\mathcal{U}}$ where $S_{\mathcal{U}} = \sum_{b \in \mathcal{U}} v_b$. Under block orthogonality,
 573

$$575 \quad \mathbb{E}\|\widehat{\xi}\|^2 = \sum_{b \in \mathcal{C}} v_b^2 + \frac{S_{\mathcal{U}}^2}{K - m}. \quad (18)$$

577 When blocks are nearly orthogonal, add $2 \sum_{b < c} v_b v_c \rho_{bc}$ to capture off-diagonal corrections.

580 **D COHERENCE-BASED ROBUSTNESS BOUND**
 581

582 We quantify deviations from the independent-marginal term in the variance due to non-orthogonality
 583 and dependence. Using the variance identity and triangle inequality,

$$585 \quad \left| \text{Var}(\widehat{\xi}) - \sum_b \left(\frac{1}{\pi_b} - 1\right) v_b^2 \right| = \left| 2 \sum_{b < c} \left(\frac{\pi_{bc}}{\pi_b \pi_c} - 1\right) v_b v_c \rho_{bc} \right| \quad (19)$$

$$588 \quad \leq 2 \sum_{b < c} \left| \frac{\pi_{bc}}{\pi_b \pi_c} - 1 \right| v_b v_c |\rho_{bc}| \quad (20)$$

$$590 \quad \leq 2\mu \sum_{b < c} \left| \frac{\pi_{bc}}{\pi_b \pi_c} - 1 \right| v_b v_c. \quad (21)$$

593 Under Bernoulli, the factor is zero because $\pi_{bc} = \pi_b \pi_c$. Under exact- K with negative association,
 the factor is bounded by the strength of repulsion and the magnitudes of the off-diagonal cosines.

594 E LOWER BOUND PROOF
595596 Starting from the variance identity,
597

598
$$\mathbb{E}\|\hat{\xi}\|^2 = \|\xi\|^2 + \sum_b \left(\frac{1}{\pi_b} - 1\right) v_b^2 + 2 \sum_{b < c} \left(\frac{\pi_{bc}}{\pi_b \pi_c} - 1\right) v_b v_c \rho_{bc}. \quad (22)$$

599

600 Lower-bound the last term by $-2 \sum_{b < c} v_b v_c$ and obtain
601

602
$$\mathbb{E}\|\hat{\xi}\|^2 \geq \|\xi\|^2 + \sum_b \frac{v_b^2}{\pi_b} - \sum_b v_b^2 - 2 \sum_{b < c} v_b v_c = \|\xi\|^2 + \sum_b \frac{v_b^2}{\pi_b} - \left(\sum_b v_b\right)^2. \quad (23)$$

603

604 By Cauchy–Schwarz, $\sum_b v_b^2 / \pi_b \geq (\sum_b v_b)^2 / \sum_b \pi_b = (\sum_b v_b)^2 / K$. Substituting proves the claim.
605606 F EXACT- K VARIANCE COMPARISON PROOF
607608 Let Var_{Bern} denote the variance under Bernoulli sampling with marginals π_b . Then $\text{Var}_{\text{Bern}} =$
609 $\sum_b \left(\frac{1}{\pi_b} - 1\right) v_b^2$. Under projection k -DPP,
610

611
$$\text{Var}_{\text{DPP}} = \sum_b \left(\frac{1}{\pi_b} - 1\right) v_b^2 - 2 \sum_{b < c} \frac{(K_{bc}^\sharp)^2}{\pi_b \pi_c} v_b v_c \rho_{bc}. \quad (24)$$

612

613 The difference $\text{Var}_{\text{Bern}} - \text{Var}_{\text{DPP}}$ equals the signed sum on the right. Under the aggregate alignment
614 assumption, this difference is nonnegative, and it is strictly positive if there exists a pair with positive
615 cosine and nonzero kernel entry.
616617 G WALL-CLOCK RULE DERIVATION AND SCHEDULING
618619 We model the per-step cost as $C_0 + C_1 K$. Under independent water-filling, write $G(K) \approx \bar{a} + \bar{b}/K$
620 for slowly varying averages \bar{a} and \bar{b} . The time to a target tolerance behaves like $\mathcal{T}(K) \propto (C_0 +$
621 $C_1 K)(\bar{a} + \bar{b}/K)$. Differentiating and setting to zero yields
622

623
$$(C_0 + C_1 K) \left(-\frac{\bar{b}}{K^2}\right) + C_1 \left(\bar{a} + \frac{\bar{b}}{K}\right) = 0 \quad \Rightarrow \quad K^* = \sqrt{\frac{C_0 \bar{b}}{C_1 \bar{a}}}. \quad (25)$$

624

625 We adopt a myopic schedule that updates K_t using exponential moving averages of $a_t = \sum_b v_{t,b}^2$,
626 and $b_t = (\sum_b v_{t,b})^2$, clips K_t to $[1, B]$, and refreshes the averages periodically to adapt to changing
627 regimes. Under exact- K , adjust the cost to include selection overhead in C_0 and, when alignment
628 holds, replace \bar{b} by an empirically reduced effective quantity.
629630 H ALGORITHMS
631632 This section provides pseudocode for the independent design with water-filling and for projection
633 k -DPP sampling.
634635 **Algorithm 1** GeoBIS one step with water-filling and probability floor
636

637 1: Input: current point X , blocks V_b , desired budget K , floor p_{\min}
638 2: Compute block gradients $\xi_{(b)}$ and scores $v_b = \|\xi_{(b)}\|_{g_X}$
639 3: Sort v_b in descending order
640 4: Find capped set \mathcal{C} and threshold λ such that $p_b^* = \min\{1, \lambda v_b\}$ satisfy $\sum_b p_b^* = K$
641 5: Set $p_b = \min\{1, \max\{p_{\min}, p_b^*\}\}$ for all b
642 6: Sample $I_b \sim \text{Bernoulli}(p_b)$ independently
643 7: Form $\hat{\xi} = \sum_b (I_b / p_b) \xi_{(b)}$
644 8: Update $X \leftarrow \text{Retr}_X(-\eta \hat{\xi})$

648 **Algorithm 2** Projection k -DPP sampling with a candidate pool

 649 1: Input: current X , blocks V_b , target size K , candidate multiplier $\alpha \geq 1$
 650 2: Compute scores v_b and unit directions $u_b = \xi_{(b)}/v_b$ for nonzero v_b
 651 3: Form a candidate set of size αK with the largest v_b
 652 4: Estimate cosines $\rho_{bc} = \langle u_b, u_c \rangle_{g_X}$ within the candidate set
 653 5: Form $M = WGW$ with $W = \text{diag}(\sqrt{v_b})$ over the candidate set
 654 6: Compute the top- K eigenspace V_K of M (randomized SVD or Oja updates)
 655 7: Set $K^\sharp = V_K V_K^\top$ and sample a projection k -DPP subset S
 656 8: Use $\pi_b = K_{bb}^\sharp$ and form $\hat{\xi} = \sum_{b \in S} \xi_{(b)}/\pi_b$
 657 9: Update $X \leftarrow \text{Retr}_X(-\eta \hat{\xi})$

660 I ORTHOGONAL AND STIEFEL DETAILS

 661
 662 On $O(d)$, the Riemannian gradient can be written as ΩX with Ω skew-symmetric. Disjoint pair
 663 blocks $E_{ij} - E_{ji}$ lead to local scores $v_{ij} = \sqrt{2}|\Omega_{ij}|$. Cayley-based retractions on small skew
 664 subspaces are efficient and parallelizable. On $\text{St}(d, p)$, the projection of the Euclidean gradient is
 665 $G - X \text{sym}(X^\top G)$, and row or column blocks are aligned with the metric. QR and polar retractions
 666 are standard choices. In both cases, scoring is local and inexpensive compared to backpropagation,
 667 which explains why independent water-filling is a practical default under short and medium sequence
 668 lengths.
 669
 670
 671
 672
 673
 674
 675
 676
 677
 678
 679
 680
 681
 682
 683
 684
 685
 686
 687
 688
 689
 690
 691
 692
 693
 694
 695
 696
 697
 698
 699
 700
 701