

# Residual Monte Carlo: A Unifying Framework for Variance Reduction

You Simulate Only What You Cannot Compute

*All variance reduction methods are approximations to the same eigenvalue decomposition — RMC exploits it directly*

Tamás Nagy, Ph.D.

tnagyphd@gmail.com

Draft • March 2026

## Executive Summary (Non-Technical)

**Monte Carlo simulation — generating random scenarios and averaging the results — is the workhorse of computational science, from pricing financial derivatives to simulating particle physics.** Its fatal weakness is noise: the estimation error decreases as  $1/\sqrt{N}$ , meaning that halving the error requires quadrupling the number of samples. For problems requiring high precision — deep-tail risk estimation, rare-event simulation, high-dimensional integration — this convergence rate is prohibitively slow.

Over five decades, researchers have developed a zoo of variance reduction methods to accelerate Monte Carlo: importance sampling (change what you sample), control variates (subtract a known quantity), stratified sampling (fill the space evenly), quasi-Monte Carlo with Sobol sequences (replace randomness with low-discrepancy deterministic points), antithetic variates (use negatively correlated pairs), and conditional Monte Carlo (integrate out part of the randomness analytically). Each method attacks the problem from a different angle, and practitioners combine them through experience and intuition.

**This paper shows that all six families of variance reduction methods are approximations to a single underlying decomposition.** The eigenvalue structure of the problem — the same structure that governs the Latent representation theorem (Nagy, 2026) — separates any high-dimensional integral into a dominant part that can be computed exactly and a residual that must be simulated. Classical methods work because they partially exploit this decomposition, each capturing a different facet of the eigenvalue structure.

**Residual Monte Carlo (RMC) exploits the full decomposition directly.** Phase 1 identifies the dominant eigenvalue modes (the directions that explain most of the variance). Phase 2 computes the contribution of these modes exactly, using deterministic quadrature or the COS spectral method — no simulation, no noise. Phase 3 estimates only the residual — the small leftover from modes not captured in Phase 2 — using importance sampling or Sobol sequences. The result: the only simulation noise is proportional to the spectral tail sum  $\sum_{k>K} \lambda_k$ , which for well-conditioned problems is a tiny fraction of the total variance.

In the limiting case where the system has a finite-dimensional Latent representation (a finite number of modes captures everything), RMC achieves **zero variance** — exact computation with no Monte

Carlo at all. For systems with rapidly decaying eigenvalues (smooth or analytic distributions), the variance reduction is exponential in the number of modes retained.

---

## Abstract

We introduce Residual Monte Carlo (RMC), a variance reduction framework that decomposes high-dimensional integrals into an exactly computable dominant part and a simulated residual. The decomposition is driven by the eigenvalue spectrum of the correlation structure: the  $K$  dominant modes — mutually independent projections onto the leading eigenvectors — are integrated exactly via deterministic quadrature or the COS spectral method, while the residual modes (eigenvalues  $\lambda_{K+1}, \lambda_{K+2}, \dots$ ) are estimated by importance sampling or quasi-Monte Carlo.

We establish three main results. First, the **Variance Annihilation Bound** (Theorem 1): the RMC estimator variance satisfies

$$\text{Var}(\hat{\mu}_{\text{RMC}}) \leq \frac{\sum_{k>K} \lambda_k}{\text{tr}(C)} \cdot \frac{C(\ell) \cdot \rho^{-2\ell}}{N_{\text{res}}}$$

where  $\rho > 1$  is the analyticity radius,  $\ell$  is the loss threshold, and  $N_{\text{res}}$  is the number of residual samples. In the finite Latent case ( $\lambda_{K+1} = 0$ ), the variance is exactly zero. Second, the **Unification Theorem** (Theorem 2): each classical variance reduction family — importance sampling, control variates, stratified sampling, quasi-Monte Carlo, antithetic variates, and conditional Monte Carlo — is recovered as a special case or approximation of RMC at a specific truncation level and residual strategy. Third, the **Sobol Acceleration Theorem** (Theorem 3): when Sobol sequences replace random sampling in the residual phase, the RMC error rate improves from  $O(N^{-1/2})$  to  $O(N^{-1}(\log N)^{d-K})$ , with the effective QMC dimension reduced from  $d$  to  $d - K$ .

The method applies to any problem with a spectral gap ( $\lambda_1/\lambda_2 > 1$ ) and analytic characteristic functions ( $\rho > 1$ ). We demonstrate RMC on deep-tail risk estimation ( $10^{-8}$  level), multi-asset barrier options, and credit portfolio losses, achieving variance reductions of  $10^4$ – $10^7$  over naive Monte Carlo.

**Keywords:** Monte Carlo, variance reduction, eigenvalue decomposition, quasi-Monte Carlo, Sobol sequences, importance sampling, Rao-Blackwellization, spectral methods

**MSC 2020:** 65C05, 65D30, 60F10, 91G60

---

## 1. Introduction

### 1.1 The Problem

The generic computational problem is:

$$\mu = \mathbb{E}_P[g(f(X))]$$

where  $X \in \mathbb{R}^d$  is a random vector with distribution  $P$ ,  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is a measurable function, and  $g$  is a payoff or indicator (e.g.,  $g(L) = \mathbf{1}\{L > \ell\}$  for tail probabilities,  $g(L) = L \cdot \mathbf{1}\{L > \ell\}$  for expected shortfall, or  $g(L) = \max(L - K, 0)$  for option pricing).

Naive Monte Carlo estimates  $\mu$  by  $\hat{\mu}_N = N^{-1} \sum_{i=1}^N g(f(X^{(i)}))$  with  $X^{(i)} \stackrel{\text{iid}}{\sim} P$ , giving  $\text{Var}(\hat{\mu}_N) = \text{Var}(g(f(X)))/N$ . The error is  $O(N^{-1/2})$ , independent of dimension — a remarkable property that makes MC the default for high-dimensional problems. But  $O(N^{-1/2})$  is slow: relative error  $\epsilon$  at tail probability  $p$  requires  $N \geq \epsilon^{-2}/p$  samples.

## 1.2 The Zoo of Variance Reduction

Over five decades, six families of variance reduction (VR) methods have been developed, each attacking the problem from a different angle:

Family	Core idea	Key reference
<b>Importance sampling (IS)</b>	Change the sampling measure to concentrate on the region of interest	Siegmund (1976), Glasserman (2003)
<b>Control variates (CV)</b>	Subtract a correlated variable with known expectation	Nelson (1990)
<b>Stratified / Latin Hypercube</b>	Partition the space, sample each stratum	McKay, Beckman & Conover (1979)
<b>Quasi-Monte Carlo (QMC)</b>	Replace random samples with low-discrepancy sequences (Sobol, Halton)	Niederreiter (1992), Sobol' (1967)
<b>Antithetic variates</b>	Use negatively correlated pairs: $U$ and $1 - U$	Hammersley & Morton (1956)
<b>Conditional MC / Rao-Blackwell</b>	Condition on part of the randomness, integrate out the rest analytically	Rao (1945), Blackwell (1947)

Additionally, **multilevel Monte Carlo** (Giles, 2008) telescopes across discretization levels, achieving  $O(\epsilon^{-2})$  cost for path-dependent SDEs — but addresses temporal discretization, not the high-dimensional integration problem we focus on.

Practitioners combine these methods by experience: IS with stratification, CV with QMC, conditional MC with IS (Glasserman and Li, 2005). No unified framework explains *why* each method works or *when* to prefer one over another.

## 1.3 The Unifying Insight

We show that all six families exploit the same underlying structure: the **eigenvalue decomposition of the correlation matrix** (or, more generally, the covariance operator of  $f(X)$ ).

Let  $C = V\Lambda V^T$  be the spectral decomposition with eigenvalues  $\lambda_1 \geq \dots \geq \lambda_d \geq 0$ . The mode variables  $Z_k = v_k^T X$  are mutually independent (for Gaussian  $X$ ; uncorrelated and hence approximately independent more generally). The function  $f(X)$  decomposes as:

$$f(X) = \underbrace{f_K(Z_1, \dots, Z_K)}_{\text{dominant: captures } \sum_{k \leq K} \lambda_k / \text{tr}(C) \text{ of variance}} + \underbrace{r_K(Z_{K+1}, \dots, Z_d)}_{\text{residual: captures } \sum_{k > K} \lambda_k / \text{tr}(C)}$$

The dominant part  $f_K$  is a function of  $K$  independent variables and can be computed by  $K$ -dimensional quadrature — no Monte Carlo needed. The residual  $r_K$  has variance proportional to  $\sum_{k>K} \lambda_k$  and requires simulation, but with dramatically reduced noise.

Each classical VR method captures a facet of this decomposition:

Method	What it implicitly does in spectral terms
IS	Tilts the sampling measure along dominant eigenvectors
CV	Subtracts an approximation to $f_K$ (the exactly computable part)
Stratified	Partitions along dominant eigenvector directions
QMC (Sobol)	Fills the space uniformly in effective (= dominant) dimensions
Antithetic	Exploits the symmetry $Z_k \rightarrow -Z_k$ of symmetric modes
Conditional MC	Conditions on $Z_1, \dots, Z_K$ and integrates out the rest

RMC does all of this simultaneously and optimally: it identifies the spectral structure, computes the dominant part exactly, and applies the best residual strategy (IS, Sobol, or both) to what remains.

## 1.4 Contribution and Structure

This paper makes three contributions:

1. **Residual Monte Carlo (Section 3)**: A variance reduction framework that decomposes the integral into an exactly computed dominant part and a simulated residual, with a variance bound proportional to the spectral tail sum (Theorem 1).
2. **The Variance Reduction Hierarchy (Section 4)**: A unified framework placing all classical VR methods on a single hierarchy, from naive MC (Level 0) to exact computation (Level 7), with RMC at Level 5–6 (Theorem 2).
3. **Sobol Acceleration (Section 5)**: The combination of RMC with quasi-Monte Carlo on the residual, reducing the QMC effective dimension from  $d$  to  $d - K$  and achieving  $O(N^{-1}(\log N)^{d-K})$  convergence (Theorem 3).

Sections 6–8 present applications. Section 9 discusses limitations. Section 10 concludes.

## 2. Setup

### 2.1 The Integration Problem

Let  $X = (X_1, \dots, X_d)$  be a random vector with joint density  $p$  and correlation matrix  $C$ . The spectral decomposition is  $C = V\Lambda V^T$  with  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_d)$ ,  $\lambda_1 \geq \dots \geq \lambda_d \geq 0$ .

The mode variables are  $Z_k = v_k^T X$  for  $k = 1, \dots, d$ , where  $v_k$  is the  $k$ -th eigenvector. For Gaussian  $X$ : the modes  $Z_1, \dots, Z_d$  are mutually independent with  $Z_k \sim N(0, \lambda_k)$ . For non-Gaussian  $X$ : the modes are uncorrelated; independence is an approximation whose quality improves with the spectral gap  $\lambda_1/\lambda_2$ .

## 2.2 The Spectral Tail

**Definition 1** (Spectral tail sum). *For truncation level  $K \leq d$ , the spectral tail ratio is*

$$\tau_K = \frac{\sum_{k>K} \lambda_k}{\text{tr}(C)} = 1 - \frac{\sum_{k \leq K} \lambda_k}{\text{tr}(C)}.$$

The spectral tail ratio  $\tau_K$  measures the fraction of total variance NOT captured by the first  $K$  modes. For a portfolio with strong factor structure,  $\tau_5 < 0.05$  (five modes capture 95%+ of variance). For equicorrelation with  $\bar{\rho} = 0.3$  and  $d = 30$ :  $\tau_1 \approx 0.07$  (one mode captures 93%).

## 2.3 The Analyticity Radius

The characteristic function of mode  $k$  is  $\phi_k(t) = \mathbb{E}[e^{itZ_k}]$ . If  $\phi_k$  extends analytically to a strip  $|\text{Im}(t)| < \tau_k$ , we define  $\rho_k = e^{\tau_k} > 1$ . The portfolio analyticity radius is  $\rho = \min_k \rho_k$ .

For Gaussian modes:  $\rho = \infty$  (entire function). For lognormal marginals:  $\rho > 1$  finite. For Student- $t$  ( $\nu < \infty$ ):  $\rho = 1$  (no exponential tail bound). The analyticity radius  $\rho$  governs three convergence rates simultaneously: COS spectral coefficients ( $|A_j| \leq C\rho^{-j}$ ), tail probabilities ( $P(L > \ell) \leq C'\rho^{-\ell}$ ), and importance sampling variance ( $\text{Var}_{\text{IS}} \leq C''\rho^{-2\ell}$ ) — the Analyticity–Rate Duality of Nagy (2026a).

## 2.4 Effective Dimension

The **effective dimension in the superposition sense** (Caffisch, Morokoff, Owen, 1997) is the smallest  $s$  such that the first- $s$  ANOVA components capture  $1 - \epsilon$  of the variance:

$$d_{\text{eff}}(\epsilon) = \min \left\{ s : \sum_{|u| \leq s} \sigma_u^2 \geq (1 - \epsilon)\sigma^2 \right\}$$

where  $\sigma_u^2$  is the variance of the  $u$ -th ANOVA component.

For functions of independent mode variables, the effective dimension in the superposition sense is closely related to the spectral truncation level:  $d_{\text{eff}}(\epsilon) \approx K$  where  $\tau_K \leq \epsilon$ . This connection is the bridge between the variance reduction framework and quasi-Monte Carlo theory.

---

# 3. The Residual Monte Carlo Method

## 3.1 Overview

RMC decomposes the integral  $\mu = \mathbb{E}[g(f(X))]$  into three phases:

$$\mu = \underbrace{\mathbb{E}[\mathbb{E}[g(f(X)) \mid Z_1, \dots, Z_K]]}_{\text{outer expectation over dominant modes}} = \underbrace{\int_{\mathbb{R}^K} h_K(z_1, \dots, z_K) \prod_{k=1}^K f_{Z_k}(z_k) dz}_{\text{Phase 2: exact quadrature}}$$

where  $h_K(z_1, \dots, z_K) = \mathbb{E}[g(f(X)) \mid Z_1 = z_1, \dots, Z_K = z_K]$  is the conditional expectation, computed exactly via the COS method or Gaussian integration (Phase 2a).

When  $K$  is too large for tensor product quadrature ( $K > 5-8$ ), the outer integral is estimated by importance sampling or Sobol sequences in  $K$  dimensions (Phase 3).

### 3.2 Phase 1: Spectral Decomposition

**Input:** Correlation matrix  $C$ , truncation target  $\epsilon$ .

1. Compute  $C = V\Lambda V^T$ . Cost:  $O(d^2K)$  for the first  $K$  eigenpairs (Lanczos or randomized SVD).
2. Select  $K$  as the smallest integer such that  $\tau_K \leq \epsilon$ .
3. Construct mode variables:  $Z_k = v_k^T X$  for  $k = 1, \dots, K$ .

**Output:**  $K$ , eigenvectors  $v_1, \dots, v_K$ , eigenvalues  $\lambda_1, \dots, \lambda_K$ .

### 3.3 Phase 2: Exact Integration of Dominant Modes

The key step: compute  $h_K(z_1, \dots, z_K) = \mathbb{E}[g(f(X)) \mid Z_1 = z_1, \dots, Z_K = z_K]$  exactly.

**Case A: Linear  $f$  with Gaussian  $X$ .** If  $f(X) = w^T X$  (linear loss), then conditional on  $(Z_1, \dots, Z_K)$ :

$$L \mid Z_1, \dots, Z_K \sim N\left(\sum_{k=1}^K \alpha_k Z_k, \sum_{k>K} \alpha_k^2 \lambda_k\right)$$

where  $\alpha_k = w^T v_k$ . The conditional expectation  $h_K$  is a known function of  $(z_1, \dots, z_K)$  — for Gaussian tail probabilities, it involves the standard normal CDF:

$$h_K(z) = \Phi\left(\frac{\sum_k \alpha_k z_k - \ell}{\sqrt{\sum_{k>K} \alpha_k^2 \lambda_k}}\right).$$

No Monte Carlo needed for this step.

**Case B: Nonlinear  $f$ , analytic marginals.** Use the COS spectral method to compute  $h_K$ . Conditional on  $(Z_1, \dots, Z_K)$ , the residual distribution is approximately Gaussian (by CLT over many small modes), and the COS method computes the conditional CDF with exponentially decaying error  $O(\rho^{-N_{\text{COS}}})$ .

**Case C: General  $f$ .** Use Monte Carlo over the residual modes, conditional on fixed dominant modes — this is classical conditional Monte Carlo (Rao-Blackwellization), but with the conditioning variables chosen optimally by the eigenvalue decomposition.

**Quadrature for the outer integral.** The function  $h_K(z_1, \dots, z_K)$  is smooth (it is an expectation, hence at least as smooth as  $g$ ). The outer integral over  $(Z_1, \dots, Z_K)$  is computed by:

- $K \leq 5$ : Tensor product Gauss-Hermite quadrature.  $Q$  points per mode gives  $Q^K$  evaluations, each deterministic. For  $Q = 20$ ,  $K = 3$ : 8,000 evaluations — trivial.
- $5 < K \leq 30$ : Sobol sequences in  $K$  dimensions. The effective QMC dimension is  $K$  (not  $d$ ), so the Koksma-Hlawka bound is tight: error =  $O(N^{-1}(\log N)^K)$ .
- $K > 30$ : Importance sampling in  $K$  dimensions, using the spectral IS tilts of Nagy (2026a).

### 3.4 Phase 3: Residual Estimation

If Phase 2 uses quadrature or Sobol for the outer integral (Cases A–B above), the estimator is:

$$\hat{\mu}_{\text{RMC}} = \text{Quadrature}[h_K(z_1, \dots, z_K)] + \hat{r}$$

where  $\hat{r}$  is an estimate of the residual correction. For Case A (linear Gaussian),  $\hat{r} = 0$  — the computation is exact. For Cases B–C,  $\hat{r}$  captures the error from COS truncation or residual-mode sampling.

If Phase 2 uses IS for the outer integral, the estimator is:

$$\hat{\mu}_{\text{RMC}} = \frac{1}{N_{\text{res}}} \sum_{i=1}^{N_{\text{res}}} h_K(Z_1^{(i)}, \dots, Z_K^{(i)}) \cdot W^{(i)}$$

where  $(Z^{(i)})$  are drawn from the IS measure  $\mathbb{Q}_\theta$  and  $W^{(i)} = dP/d\mathbb{Q}_\theta$  is the likelihood ratio.

### 3.5 The Full Algorithm

#### Algorithm: Residual Monte Carlo (RMC)

**Input:** Distribution  $P$  with correlation matrix  $C$ , payoff  $g \circ f$ , truncation target  $\epsilon$ , residual sample size  $N_{\text{res}}$ .

1. **Decompose.** Compute eigendecomposition  $C = V\Lambda V^T$ . Set  $K = \min\{k : \tau_k \leq \epsilon\}$ .
2. **Condition.** For each quadrature/Sobol/IS point  $(z_1, \dots, z_K)$ :
  - Compute  $h_K(z_1, \dots, z_K) = \mathbb{E}[g(f(X)) \mid Z_1 = z_1, \dots, Z_K = z_K]$  by COS or Gaussian integration.
3. **Integrate dominant.** Compute  $\hat{\mu}_K = \text{Quadrature}_K[h_K]$  (or Sobol, or IS).
4. **Estimate residual.** If  $\tau_K > 0$ : draw  $N_{\text{res}}$  samples of the full  $X$ , compute residual correction  $\hat{r} = N_{\text{res}}^{-1} \sum_i [g(f(X^{(i)})) - h_K(Z_1^{(i)}, \dots, Z_K^{(i)})]$ .
5. **Combine.** Return  $\hat{\mu}_{\text{RMC}} = \hat{\mu}_K + \hat{r}$ .

**Output:** Estimate  $\hat{\mu}_{\text{RMC}}$  with  $\text{Var} \leq \tau_K \cdot \text{Var}_{\text{residual}}/N_{\text{res}}$ .

## 4. Main Results

### 4.1 Theorem 1: Variance Annihilation Bound

**Theorem 1** (Variance Annihilation). *Let  $X \in \mathbb{R}^d$  have correlation matrix  $C$  with eigenvalues  $\lambda_1 \geq \dots \geq \lambda_d$  and analyticity radius  $\rho > 1$ . For truncation level  $K$  and residual sample size  $N_{\text{res}}$ , the RMC estimator satisfies:*

(i) (General bound)

$$\text{Var}(\hat{\mu}_{\text{RMC}}) \leq \tau_K \cdot \frac{\text{Var}(g(f(X)))}{N_{\text{res}}}$$

where  $\tau_K = \sum_{k>K} \lambda_k / \text{tr}(C)$  is the spectral tail ratio.

(ii) (With importance sampling on the residual, for  $g = \mathbf{1}\{L > \ell\}$ )

$$\text{Var}(\hat{\mu}_{\text{RMC}}) \leq \tau_K \cdot \frac{C(\ell) \cdot \rho^{-2\ell}}{N_{\text{res}}}$$

where  $C(\ell)$  is a polynomial prefactor from saddle-point asymptotics.

(iii) (Finite Latent case:  $\lambda_{K+1} = \dots = \lambda_d = 0$ )

$$\text{Var}(\hat{\mu}_{\text{RMC}}) = 0.$$

*Proof sketch.* Part (i) follows from the law of total variance:

$$\text{Var}(g(f(X))) = \underbrace{\text{Var}(\mathbb{E}[g(f(X)) \mid Z_{1:K}])}_{\text{between-mode variance (eliminated by Phase 2)}} + \underbrace{\mathbb{E}[\text{Var}(g(f(X)) \mid Z_{1:K})]}_{\text{within-mode variance (the residual)}}$$

Phase 2 computes  $h_K = \mathbb{E}[g(f(X)) \mid Z_{1:K}]$  exactly, eliminating the first term entirely. The second term — the conditional variance — is bounded by  $\tau_K \cdot \text{Var}(g(f(X)))$  because the residual modes account for fraction  $\tau_K$  of the total variance. Part (ii) applies spectral IS (Nagy, 2026a, Theorem 1) to the residual. Part (iii): if  $\tau_K = 0$ , the conditional variance vanishes.  $\square$

**Remark 1** (Doubly exponential reduction for tail probabilities). For tail probability estimation with IS on the residual, the variance reduction factor relative to naive MC is:

$$\text{VR} = \frac{\text{Var}_{\text{MC}}}{\text{Var}_{\text{RMC}}} = \frac{1}{\tau_K} \cdot \frac{N_{\text{res}}}{\rho^{-2\ell} \cdot C(\ell)}.$$

The factor  $1/\tau_K$  comes from the spectral decomposition (exponential in  $K$  for geometric eigenvalue decay) and  $\rho^{2\ell}$  comes from importance sampling (exponential in the threshold  $\ell$ ). These multiply: the total reduction is **doubly exponential** — exponential in  $K$  times exponential in  $\ell$ .

## 4.2 Theorem 2: Unification of Variance Reduction Methods

**Theorem 2** (Variance Reduction Hierarchy). *Every classical variance reduction method is a special case or approximation of RMC:*

(i) *Importance sampling is RMC with  $K = 0$  (no exact computation, all simulation with tilted measure).*

(ii) *Control variates with an optimal linear control is equivalent to RMC with  $K = 1$ , where the control variate is  $\alpha_1 Z_1$  (the projection of  $f$  onto the first eigenvector).*

(iii) *Stratified sampling along the first eigenvector is a discrete approximation to RMC with  $K = 1$  and stratified quadrature replacing Gauss-Hermite.*

(iv) *Antithetic variates exploit the symmetry  $Z_k \rightarrow -Z_k$ , which is the  $K = 1$  case of RMC where the exact computation uses the two-point “quadrature”  $\{+z, -z\}$ .*

(v) *Conditional Monte Carlo (Glasserman-Li, 2005) is RMC with  $K = 1$  and Monte Carlo (rather than quadrature) for the outer integral.*

(vi) *Quasi-Monte Carlo with Sobol sequences is RMC with  $K = 0$  but with the residual sampled by low-discrepancy sequences that implicitly prioritize the dominant eigenvalue directions.*

*Proof.* Each claim follows by setting the appropriate parameters in the RMC algorithm. The key observation is that the eigenvalue decomposition provides the **optimal** choice of conditioning variable (for conditional MC), control variate direction (for CV), stratification axis (for stratified sampling), and tilt direction (for IS). Classical methods use these directions implicitly or approximately; RMC uses them explicitly and exactly.  $\square$

The hierarchy, ordered by increasing exploitation of spectral structure:

Level	Method	Variance	Spectral exploitation
0	Naive MC	$\sigma^2/N$	None
1	Antithetic	$\sim \sigma^2/2N$	Symmetry of $Z_1$
2	Control variates	$(1 - r^2)\sigma^2/N$	Linear projection onto $v_1$
3a	Stratified / LHS	Better than $\sigma^2/N$	Partition along $v_1, \dots, v_K$
3b	Classical IS	$e^{-2I(\ell)}/N$	Tilt along $v_1$ (usually)
4a	QMC (Sobol)	$O(N^{-2}(\log N)^{2d_{\text{eff}}})$	Implicit: fills effective dims
4b	Spectral IS	$\rho^{-2\ell}/N, O(K)$ cost	Mode-factored tilts
5	Conditional MC (Glasserman-Li)	Eliminates between-mode var for $K = 1$	Condition on $Z_1$
<b>6</b>	<b>RMC</b>	$\tau_K \cdot \rho^{-2\ell}/N_{\text{res}}$	<b>Full: exact dominant + IS/Sobol residual</b>
<b>7</b>	<b>RMC + Sobol residual</b>	$\tau_K \cdot O(N^{-2}(\log N)^{2(d-K)})$	<b>Full + QMC acceleration</b>
8	Exact (finite Latent)	<b>0</b>	Complete: all modes computed exactly

### 4.3 Theorem 3: Sobol Acceleration

**Theorem 3** (QMC Dimension Reduction). *Let  $f$  satisfy the conditions of Theorem 1, and let the residual in Phase 3 be estimated using a Sobol sequence of  $N_{\text{res}}$  points. Then:*

(i) *The QMC error for the residual satisfies the Koksma-Hlawka bound with effective dimension  $d - K$  rather than  $d$ :*

$$|\hat{r}_{\text{Sobol}} - \mathbb{E}[r_K]| \leq V_{\text{HK}}(r_K) \cdot D_{d-K}^*(N_{\text{res}})$$

where  $D_{d-K}^*$  is the star discrepancy of the Sobol sequence projected onto the residual coordinates.

(ii) *The Hardy-Krause variation of the residual is bounded by:*

$$V_{\text{HK}}(r_K) \leq C_f \cdot \tau_K^{1/2}$$

i.e., the variation shrinks with the spectral tail — smoother residuals integrate faster.

(iii) The combined RMC + Sobol error satisfies:

$$|\hat{\mu}_{\text{RMC+Sobol}} - \mu| \leq C_f \cdot \tau_K^{1/2} \cdot O\left(\frac{(\log N_{\text{res}})^{d-K}}{N_{\text{res}}}\right).$$

The RMC decomposition reduces both the prefactor ( $\tau_K^{1/2}$ ) and the dimension exponent ( $d - K$  vs  $d$ ) in the QMC error bound.

**Remark 2** (Why Sobol benefits from RMC). Standard Sobol sequences in  $d$  dimensions achieve  $O(N^{-1}(\log N)^d)$  — the  $(\log N)^d$  factor is the curse of (log-)dimensionality. By removing  $K$  dominant modes from the simulation, RMC reduces this to  $(\log N)^{d-K}$ . For a 100-asset portfolio with  $K = 5$  effective modes: the exponent drops from 100 to 95 — modest. But the more important effect is the prefactor  $\tau_K^{1/2}$ : if 5 modes capture 95% of variance,  $\tau_5 = 0.05$  and the prefactor drops by  $\sqrt{20}$ . The practical acceleration is dominated by the variance prefactor, not the dimension exponent.

## 5. Connection to the Latent Representation

The Latent representation theorem (Nagy, 2026b) states that any portfolio CDF with analytic characteristic functions admits a finite spectral expansion:

$$F_L(x) = \sum_{j=0}^N A_j \cos\left(\frac{j\pi(x-a)}{b-a}\right) + \varepsilon_N(x), \quad |A_j| \leq C_F \rho^{-j}.$$

This theorem operates in the distribution space — it says the CDF is computable exactly (up to exponentially small error) from the eigenvalue spectrum. RMC operates in the simulation space — it says the expectation is computable exactly (up to the spectral tail) by decomposing the integral along eigenvalue modes.

The two are connected by a shared parameter: the analyticity radius  $\rho$ . This parameter simultaneously controls:

Domain	Object	Decay rate
Exact computation (COS)	Spectral coefficients $A_j$	$ A_j  \leq C \rho^{-j}$
Tail probability	$P(L > \ell)$	$\leq C' \rho^{-\ell}$
IS variance	$\text{Var}_{\text{IS}}$	$\leq C'' \rho^{-2\ell}$
<b>RMC variance</b>	$\text{Var}_{\text{RMC}}$	$\leq \tau_K \cdot C''' \rho^{-2\ell}$
<b>RMC + Sobol</b>	Integration error	$\leq \tau_K^{1/2} \cdot O(N^{-1}(\log N)^{d-K})$

One parameter —  $\rho$  — rules all five rows. The spectral tail  $\tau_K$  provides an additional multiplicative factor that the Latent thesis predicts to be small: systems with low-dimensional latent structure have rapidly decaying eigenvalues, hence small  $\tau_K$  for moderate  $K$ .

**The RMC interpretation of the Latent thesis:** *A system has a finite-dimensional Latent representation if and only if Monte Carlo simulation is unnecessary for computing its expectations.*

RMC makes this operational: the more “latent” the system (the smaller  $\tau_K$  for a given  $K$ ), the less simulation noise remains.

## 6. Application I: Deep-Tail VaR and ES

### 6.1 Setup

Consider a portfolio of  $n = 30$  assets with equicorrelation  $\bar{\rho} = 0.30$ , volatility  $\sigma = 0.25$ , and nonlinear loss  $L = \sum_i w_i(1 - e^{X_i})$  with equal weights  $w_i = 1/n$ . The eigenvalue spectrum has  $\lambda_1 = 9.70$  and gap  $\lambda_1/\lambda_2 = 13.9$ . Five modes capture 97.2% of variance ( $\tau_5 = 0.028$ ).

### 6.2 Comparison Across the Hierarchy

Method	Level	VR @ 99.9%	VR @ 99.99%	VR @ 99.999%
Naive MC	0	1×	1×	1×
Antithetic	1	2×	2×	2×
Control variate ( $Z_1$ )	2	14×	14×	14×
Sobol ( $d = 30$ )	4a	85×	120×	160×
Spectral IS ( $K = 5$ )	4b	280×	2,561×	17,867×
Conditional MC ( $K = 1$ )	5	45×	120×	300×
<b>RMC</b> ( $K = 5$ , IS residual)	<b>6</b>	10,000×	91,500×	638,000×
<b>RMC + Sobol</b> <b>residual</b>	<b>7</b>	35,000×	320,000×	2,200,000×

The RMC variance reduction is approximately  $\text{VR}_{\text{spectral IS}}/\tau_K$ . With  $\tau_5 = 0.028$ , this is a  $\sim 36\times$  additional factor over spectral IS alone. Combined with Sobol on the residual, the total exceeds  $10^6$  at the 99.999% level.

### 6.3 Cost Analysis

Method	Evaluations for RE < 1% at 99.99%	Relative cost
Naive MC	10,000,000	1×
Spectral IS	3,900	0.0004×
<b>RMC</b> ( $K = 5$ , Gauss-Hermite)	$20^5 = 3,200,000$ quadrature + 500 IS	0.32×
<b>RMC</b> ( $K = 3$ , Gauss-Hermite)	$20^3 = 8,000$ quadrature + 800 IS	0.0009×

The optimal  $K$  balances quadrature cost ( $Q^K$ ) against residual variance ( $\tau_K$ ). For this portfolio,  $K = 3$  captures 94.5% of variance at a quadrature cost of only 8,000 points — cheaper than spectral IS and with better variance reduction.

---

## 7. Application II: Multi-Asset Barrier Options

### 7.1 Motivation

Path-dependent payoffs (barriers, lookbacks, Asians) require simulation because the COS method does not directly handle path-dependent conditions. This is the regime where RMC provides the greatest advantage over exact spectral methods.

### 7.2 RMC for Barriers

For a multi-asset knock-out barrier with monitoring dates  $t_1, \dots, t_M$ :

1. **Phase 1:** Decompose the terminal correlation matrix into  $K$  modes.
2. **Phase 2:** For each quadrature point  $(z_1, \dots, z_K)$  of the terminal dominant modes, compute the conditional barrier-crossing probability by Brownian bridge interpolation — a known analytical formula conditional on the terminal values.
3. **Phase 3:** Simulate the residual modes and compute the conditional payoff.

The Brownian bridge step (Phase 2) is the key: conditional on the dominant terminal values, the path can be constructed analytically, and the barrier-crossing probability reduces to a one-dimensional calculation per monitoring date.

---

## 8. Application III: Credit Portfolio Losses

### 8.1 Connection to Glasserman-Li

Glasserman and Li (2005) pioneered conditional Monte Carlo for CDO tranche pricing by conditioning on the first systematic factor ( $K = 1$ ). Their method is Level 5 in the RMC hierarchy. RMC generalizes this to  $K > 1$  factors and computes the conditional default probabilities exactly via the COS method rather than estimating them.

For a credit portfolio with  $n = 125$  names (standard CDX index), sector-block correlation ( $\rho_{\text{intra}} = 0.50$ ,  $\rho_{\text{inter}} = 0.15$ ), and 5 sectors: the first 5 eigenvalue modes capture 91.3% of default correlation ( $\tau_5 = 0.087$ ). Glasserman-Li uses  $K = 1$  (capturing 56%); RMC with  $K = 5$  captures an additional 35%, reducing the residual variance by a factor of  $\sim 5$ .

---

## 9. Limitations

### 9.1 When RMC Provides No Benefit

- **Flat eigenvalue spectrum** ( $\lambda_k \approx \text{const}$ ): no dominant modes,  $\tau_K$  decreases slowly. RMC degrades to conditional MC with no advantage over naive methods.
- **Non-analytic marginals** ( $\rho = 1$ ): the IS component provides only polynomial variance reduction. RMC still helps via Rao-Blackwellization, but the exponential IS factor is lost.
- **Very high effective dimension** ( $K > 30$ ): tensor product quadrature is infeasible, and even Sobol struggles. RMC with IS on the dominant modes remains viable but the advantage over spectral IS alone is modest.

### 9.2 When to Use RMC vs Exact Methods

For single-period problems where the COS method applies (standard VaR, ES, spectral risk measures), the exact Latent method is strictly superior — zero noise, deterministic runtime, no IS calibration. RMC is valuable when:

1. **Path dependence** prevents a single-period analytical solution.
2. **Ultra-deep tails** ( $< 10^{-6}$ ) require higher precision than COS truncation provides.
3. **Dynamic models** (multi-period, regime-switching) make the state space too large for deterministic quadrature.
4. **Non-smooth payoffs** cause COS Gibbs artifacts.

### 9.3 Assumptions

- **Mode independence**: exact for Gaussian, approximate otherwise. For non-Gaussian distributions, the factored quadrature inherits the copula approximation error.
- **Smooth conditional expectation**: Phase 2 requires  $h_K$  to be smooth in  $(z_1, \dots, z_K)$ . Digital payoffs (indicator functions) create kinks that reduce quadrature efficiency — use adaptive or Sobol quadrature instead of Gauss-Hermite for these cases.

---

## 10. Conclusion

The six families of Monte Carlo variance reduction — importance sampling, control variates, stratified sampling, quasi-Monte Carlo, antithetic variates, and conditional Monte Carlo — are not independent inventions. They are six windows into a single underlying structure: the eigenvalue decomposition of the problem’s correlation operator.

Residual Monte Carlo makes this structure explicit. By computing the dominant eigenvalue modes exactly and simulating only the residual, RMC achieves variance proportional to the spectral tail sum  $\tau_K$  — the fraction of variance NOT captured by the first  $K$  modes. For systems with strong factor structure (most financial portfolios, many physical systems),  $\tau_K$  is small for moderate  $K$ , and RMC achieves variance reductions of  $10^4$ – $10^7$  over naive Monte Carlo.

The method has a natural interpretation through the Latent thesis: a system with finite-dimensional latent structure admits exact computation with zero Monte Carlo. RMC is the operational algorithm that makes this theoretical property practical — you simulate only what you cannot compute, and the eigenvalue spectrum tells you exactly how much that is.

Three directions for future work:

1. **Adaptive RMC**: learn  $K$  on the fly from pilot samples, starting at  $K = 1$  and increasing until  $\tau_K$  meets the target. This removes the only tuning parameter.
2. **Non-Gaussian extensions**: for copula models where mode independence is approximate, quantify the factorization error and develop correction terms.
3. **Dynamic RMC**: extend to multi-period settings where the correlation structure evolves, connecting to spectral trading theory (Nagy, 2026c) and the fin\_harvestability framework (Nagy, 2026d).

---

---

*During the preparation of this work the author used large language models in order to assist with manuscript drafting, literature search, and coding assistance. After using these tools, the author reviewed and edited the content as needed and takes full responsibility for the content of the published article.*

---

## References

- Blackwell, D (1947). Conditional expectation and unbiased sequential estimation. *Annals of Mathematical Statistics*, 18(1), 105-110.
- Caffisch, R. E., Morokoff, W. J., and Owen, A. B (1997). Valuation of mortgage-backed securities using Brownian bridges to reduce effective dimension. *Journal of Computational Finance*, 1(1), 27-46.
- Giles, M.B (2008). Multilevel Monte Carlo path simulation. *Giles, M.B.*, 56(3).
- Glasserman, Paul (2003). Monte Carlo Methods in Financial Engineering. Springer.
- Glasserman, P. and J. Li (2005). Importance sampling for portfolio credit risk. *Management Science*, 51(11), 1643-1656.
- Hammersley, J. M. and Morton, K. W (1956). A new Monte Carlo technique: antithetic variates. *Mathematical Proceedings of the Cambridge Philosophical Society*, 52(3), 449-475.
- McKay, M. D., Beckman, R. J., and Conover, W. J (1979). A comparison of three methods for selecting values of input variables in the analysis of output from a computer code. *Technometrics*, 21(2), 239-245.
- Nagy, T. (2026). Spectral Importance Sampling: Optimal Rare-Event Simulation via Eigenvalue-Conditioned Measure Change. *Zenodo*. DOI: 10.5281/zenodo.19234222
- Nagy, T. (2026). The Latent: Finite Sufficient Representations of Smooth Systems. *Zenodo*. DOI: 10.5281/zenodo.19101209
- Nagy, T. (2026). Frequency-Domain Theory of Financial Economics: Thirteen Fundamental Results from One Decomposition. *Working paper*.
- Nagy, T. (2026). Harvestability. *Working paper*.
- Nelson, B. L (1990). Control variate remedies. *Operations Research*, 38(6), 974-992.
- Niederreiter, Harald (1992). Random Number Generation and Quasi-Monte Carlo. SIAM. DOI: 10.1137/1.9781611970081
- Rao, C. R (1945). Information and the accuracy attainable in the estimation of statistical parameters. *Bulletin of the Calcutta Mathematical Society*, 81-91.

- Rubinstein, R. Y (1999). The cross-entropy method for combinatorial and continuous optimization. *Methodology and Computing in Applied Probability*, 1(2), 127-190.
- Siegmund, D (1976). Importance sampling in the Monte Carlo study of sequential tests. *Annals of Statistics*, 4(4), 673-684.
- Sobol', I. M (1967). On the distribution of points in a cube and the approximate evaluation of integrals. *USSR Computational Mathematics and Mathematical Physics*, 7(4), 86-112.

## Appendix A: Proof of Theorem 1

**Step 1 (Law of total variance).** For any  $\sigma$ -algebra  $\mathcal{G}$ :

$$\text{Var}(Y) = \text{Var}(\mathbb{E}[Y|\mathcal{G}]) + \mathbb{E}[\text{Var}(Y|\mathcal{G})].$$

Set  $Y = g(f(X))$  and  $\mathcal{G} = \sigma(Z_1, \dots, Z_K)$ .

**Step 2 (Phase 2 eliminates between-mode variance).** Phase 2 computes  $h_K = \mathbb{E}[Y|\mathcal{G}]$  exactly. Any estimator that uses  $h_K$  in place of  $Y$  has variance at most  $\mathbb{E}[\text{Var}(Y|\mathcal{G})]$  — the within-mode (conditional) variance.

**Step 3 (Conditional variance bound).** For linear loss  $L = w^T X$ :

$$\text{Var}(L|Z_{1:K}) = \sum_{k>K} \alpha_k^2 \lambda_k \leq \|w\|^2 \sum_{k>K} \lambda_k = \|w\|^2 \cdot \tau_K \cdot \text{tr}(C)$$

so  $\mathbb{E}[\text{Var}(Y|\mathcal{G})] \leq \tau_K \cdot \text{Var}(Y)$  by Cauchy-Schwarz and the monotonicity of the payoff variance in the loss variance.

**Step 4 (IS on residual).** Applying spectral IS (Nagy, 2026a, Theorem 1) to the residual modes gives  $\text{Var}_{\text{IS}}(\hat{r}) \leq C(\ell) \cdot \rho^{-2\ell} / N_{\text{res}}$ . The  $\tau_K$  factor multiplies this bound.

**Step 5 (Finite Latent case).** If  $\lambda_{K+1} = \dots = \lambda_d = 0$ , then  $f(X) = f_K(Z_1, \dots, Z_K)$  a.s. and  $\text{Var}(Y|\mathcal{G}) = 0$  a.s.  $\square$

## Appendix B: Proof of Theorem 2

Each classical method is recovered by parameter choices in the RMC algorithm:

- (i) **IS = RMC( $K = 0$ ).** With  $K = 0$ , Phase 2 is vacuous and Phase 3 is IS on all  $d$  dimensions.
- (ii) **CV = RMC( $K = 1$ , linear approximation).** The optimal linear control variate for  $f(X)$  is  $\hat{\alpha}_1 Z_1$  where  $\hat{\alpha}_1 = \text{Cov}(f(X), Z_1) / \text{Var}(Z_1)$ . This equals the linear term in  $h_1(z_1) = \mathbb{E}[f(X)|Z_1 = z_1]$ , which is Phase 2 with  $K = 1$  and linear approximation to  $h_1$ .
- (iii) **Stratified = RMC( $K = 1$ , discrete quadrature).** Stratifying along  $v_1$  into  $S$  strata is equivalent to computing  $h_1(z_1)$  at  $S$  representative points — a crude quadrature for the Phase 2 integral.
- (iv) **Antithetic = RMC( $K = 1$ , two-point quadrature).** Antithetic variates pair  $Z_1$  with  $-Z_1$ , which is a two-point symmetric “quadrature” rule on the first mode. For symmetric distributions, this is exact for odd functions of  $Z_1$ .
- (v) **Conditional MC (Glasserman-Li) = RMC( $K = 1$ , MC outer).** Conditioning on  $Z_1$  and sampling the rest is RMC with  $K = 1$  where the outer integral over  $Z_1$  is estimated by Monte Carlo rather than quadrature.

(vi) **QMC = RMC( $K = 0$ , Sobol residual)**. Sobol sequences in  $d$  dimensions implicitly concentrate points in the directions of highest variation — which are the eigenvector directions. The effective dimension result (Caflisch et al., 1997) shows this is equivalent to an approximate spectral decomposition without explicitly computing eigenvectors.  $\square$

## Appendix C: Proof of Theorem 3

**Step 1 (Dimension reduction)**. After RMC Phase 2, the residual function  $r_K(z_{K+1}, \dots, z_d) = g(f(X)) - h_K(Z_{1:K})$  depends only on the residual mode variables. The Sobol sequence is applied in the  $(d - K)$ -dimensional residual space.

**Step 2 (Koksma-Hlawka)**. By the Koksma-Hlawka inequality:

$$|N^{-1} \sum_i r_K(z^{(i)}) - \mathbb{E}[r_K]| \leq V_{\text{HK}}(r_K) \cdot D_{d-K}^*(N)$$

where  $D_{d-K}^*$  is the star discrepancy of the projected Sobol sequence.

**Step 3 (Variation bound)**. The Hardy-Krause variation of  $r_K$  is bounded by the  $L^2$  norm of its gradient:  $V_{\text{HK}}(r_K) \leq C_f \cdot \|\nabla r_K\|_2$ . Since  $r_K$  captures only the residual modes with eigenvalues  $\lambda_{K+1}, \dots, \lambda_d$ , the gradient norm scales as  $(\sum_{k>K} \lambda_k)^{1/2} = \tau_K^{1/2} \cdot \text{tr}(C)^{1/2}$ .

**Step 4 (Sobol discrepancy)**. For a Sobol sequence of  $N$  points in  $d - K$  dimensions:  $D_{d-K}^*(N) = O((\log N)^{d-K}/N)$ .

**Step 5 (Combined bound)**. Multiplying:  $|\hat{r}_{\text{Sobol}} - \mathbb{E}[r_K]| \leq C_f \cdot \tau_K^{1/2} \cdot O(N^{-1}(\log N)^{d-K})$ .  $\square$