

```

1 !pip install mpi4py
2 !pip install POT
3 !nvidia-smi
4 !pip install -q --upgrade cupy-cuda12x
5 !pip install softimpute          # notice: no underscore
6 # ===== #
7 # CELL 1: Project Setup, Imports, Logging, Config
8 # ===== #
9 import os
10 import sys
11 import time
12 import math
13 import re
14 import gc
15 import logging
16 from pathlib import Path
17 from typing import Tuple, List, Dict, Optional, Union, Callable, Any
18 import numpy as np
19 import pandas as pd
20 import matplotlib.pyplot as plt
21 from scipy import sparse
22 from scipy.sparse.linalg import svds, LinearOperator # Import LinearOperator
23 from scipy.optimize import OptimizeResult # For line search return consistency
24 from numpy.random import default_rng, Generator
25 from sklearn.model_selection import train_test_split # For train/validation split
26 # --- Mount Google Drive ---
27 from google.colab import drive # Uncomment if using Colab
28 drive.mount("/content/drive", force_remount=True)
29 DRIVE_MOUNTED = True
30 # right after the imports
31 import logging
32 logging.disable(logging.WARNING) # hides all warnings emitted via logging
33
34 # === ADDED Block 5 (MPI) ===
35 try:
36     from mpi4py import MPI
37     COMM = MPI.COMM_WORLD
38     RANK_MPI = COMM.Get_rank()
39     SIZE_MPI = COMM.Get_size()
40     if RANK_MPI == 0: print(f"+++ MPI Detected: Running with {SIZE_MPI} processes. +++")
41 except ImportError:
42     COMM = None
43     RANK_MPI = 0
44     SIZE_MPI = 1
45     # print("+++ MPI Not Found: Running in serial mode. +++") # Less verbose
46
47 # === ADDED Block 6 === (Import for OT demo)
48 try:
49     import ot
50     OT_AVAILABLE = True
51 except ImportError:
52     OT_AVAILABLE = False
53     if RANK_MPI == 0: print("Warning: POT library not found. Skipping Barycentre demo.")
54
55 # === ADDED Block 6 (PCA) ===
56 try:
57     from sklearn.decomposition import PCA
58     PCA_AVAILABLE = True
59 except ImportError:
60     PCA_AVAILABLE = False
61     if RANK_MPI == 0: print("Warning: sklearn not found. Skipping PCA trajectory plot.")
62
63
64 # --- Logging Setup (Initialize Logger FIRST) ---
65 logging.basicConfig(
66     level=logging.INFO,
67     format="%(asctime)s [%(levelname)s] %(message)s",
68     handlers=[logging.StreamHandler(sys.stdout)],
69     force=True, # Overwrite any existing config
70 )
71 logger = logging.getLogger(__name__)
72
73 # --- Mount Drive ---
74 if RANK_MPI == 0: print("+++ Mounting Google Drive +++")
75 try:
76     # Only rank 0 should try to force remount if needed

```

```

77     drive.mount('/content/drive', force_remount=(RANK_MPI == 0))
78     if RANK_MPI == 0: print("Drive mounted.")
79     if COMM and SIZE_MPI > 1: COMM.Barrier() # Ensure drive is mounted
80 except Exception as e:
81     if RANK_MPI == 0: print(f"Error mounting drive: {e}")
82     if COMM and SIZE_MPI > 1: COMM.Abort()
83     raise
84
85 # --- Optional: Try importing CuPy for GPU acceleration ---
86 # NOTE: Efficient SoftImpute implementation below uses SciPy sparse ops,
87 # GPU acceleration would require re-implementing the LinearOperator with CuPy sparse.
88 try:
89     import cupy as cp
90     import cupyx.scipy.sparse as cpx
91     CUPY_AVAILABLE = False # Disable GPU for SoftImpute for now due to LinearOperator complexity
92     logger.warning("CuPy found, but GPU acceleration for efficient SoftImpute is NOT enabled in this version.")
93     if 'cp' not in locals(): cp = np
94     if 'cpx' not in locals(): cpx = sparse
95 except ImportError:
96     CUPY_AVAILABLE = False
97     cp = np ; cpx = sparse
98     logger.warning("CuPy not found, will run on CPU using NumPy/SciPy.")
99
100 logger.info("+++ Cell 1: Setup, Imports, Logging, Config +++")
101
102 # --- Global Config ---
103 # --- MOVIELENS 1M Configuration ---
104 DATA_DIR_STR = "/content/drive/MyDrive/ml-1m" # ADJUST PATH AS NEEDED
105 RATINGS_FILENAME = "ratings.dat"
106 VALIDATION_FRACTION = 0.2 # Hold out 20% for validation
107 # --- USE COMPLETE DATASET (FIX 1) ---
108 RATING_LIMIT = None # Load all ratings from ml-1m
109 RANK = 10 # Factorization rank (r in paper) for non-convex
110 LAM = 1e-2 # Regularization parameter  $\lambda$ 
111 LAM_SQ = LAM ** 2 #  $\lambda^2$  for non-convex model factor regularization
112 LAM_BIAS = 1e-4 # Regularization for bias terms
113 SEED = 0 # Use consistent seed from long.txt
114 # --- INCREASED ITERATIONS ---
115 N_ITERS_ALL = 20 # Iterations/epochs for ALL solvers
116 CONVEX_RANK_K = 50 # Max rank for Soft-Impute intermediate SVDs
117 SOFT_IMPUTE_TOL = 1e-4 # Convergence tolerance for Soft-Impute
118 N_ITERS_CONVEX = N_ITERS_ALL # Use same number of iterations for SoftImpute
119 # --- SVRG Params ---
120 INIT_LR_SVRG = 1e-3 # Base Learning rate for SVRG inner solver
121 SVRG_INNER_STEPS_DIVISOR = 1 # Use full inner pass
122 GRAD_CLIP_THRESHOLD = 10.0 # Max norm for SVRG gradients before update
123 RSVRG_BATCH_SIZE = 100 # Batch size for non-convex SVRG refresh step
124 # --- ALS Params ---
125 ALS_TOL = 1e-4 # Convergence tolerance for ALS based on RMSE change
126 ALS_MAX_ITER = N_ITERS_ALL # Use same iter count as others for comparison
127 # --- RGD/Accelerated Params ---
128 INIT_LR_RIEMANN = 0.5 # Initial LR for RGD/RAGD/Catalyst/DANE line search
129 LS_BETA = 0.5 # Line search reduction factor
130 LS_SIGMA = 1e-4 # Sufficient decrease parameter
131 RAGD_GAMMA = 1.0; RAGD_MU = 5.0; RAGD_BETA = 5.0
132 DANE_KAPPA = 1.0
133 KAPPA_0 = 1e-1; KAPPA_CVX = 1e-1; INNER_T = 5; INNER_S_BASE = 10; MAX_KAPPA_DOUBLINGS = 10
134 # --- Smaller Initialization Scale ---
135 INIT_SCALE_NON_CONVEX = 0.01 # Smaller scale for initial U, W
136 # --- Configuration from Proposal/long.txt ---
137 RETRACTION_NAME = "orthonormal" # Options: "orthonormal", "cayley", "projection"
138 REG_DISTANCE = "euclid" # Options: "euclid", "retraction"
139 INNER_SOLVER = "svrg" # Options: "svrg", "sarah", "spider" (for Catalyst)
140 ETA_GRAD = 1e-3 # Adaptive stopping tolerance for inner grad norm
141 ETA_DIST = 1e-4 # Adaptive stopping tolerance for inner step size
142 CATALYST_INNER_T_EPOCHS = 1 # Epochs for Alg phi_1 check budget
143 CATALYST_INNER_S_EPOCHS_BASE = 2 # Base epochs for S_k schedule
144 RSVRG_LR = 1e-3 # Step size for RSVRG/SARAH/SPIDER inner loops
145
146 # --- Derived Globals ---
147 GLOBAL_RNG = default_rng(SEED)
148 DATA_DIR = Path(DATA_DIR_STR)
149 I_r = np.eye(RANK, dtype=np.float64) # Identity matrix of size RANK
150
151 # Check Data Directory
152 if DRIVE_MOUNTED and not DATA_DIR.is_dir():
153     if RANK_MPI == 0: logger.warning(f"DATA_DIR '{DATA_DIR}' not found. Please check the path.")
154     raise ValueError(f"DATA_DIR '{DATA_DIR}' not found. Please check the path.")

```

```

154 elif not DRIVE_MOUNTED:
155     if RANK_MPI == 0: logger.warning(f"Google Drive not mounted.")
156
157 logger.info("Cell 1 initialisation complete.")
158
159 # ===== #
160 # CELL 2: Data Loading and Preprocessing (MovieLens 1M)
161 # ===== #
162 logger.info("+++ Cell 2: Loading and Processing Data (MovieLens 1M) +++")
163 # --- Manifold Operations ---
164 # --- universal 2-tuple helper for loss/grad (used by Catalyst) ---
165 def stochastic_gradient_batch(U, user_ids, N_users, N_movies, loss_args):
166     """
167     Vectorised version of `stochastic_gradient_single_user`.
168     Accumulates the (un-scaled) gradient over the provided user_ids.
169     """
170     G = np.zeros_like(U, dtype=np.float32)
171     for uid in user_ids:
172         G += stochastic_gradient_single_user(U, int(uid), N_users, N_movies, loss_args)
173     return G / max(1, len(user_ids)) # average over the batch
174
175 def loss_and_grad_corrected(U, W, bu, bi, *rest):
176     """ Wrapper for loss_and_grad_serial_with_biases to return only loss and grad_U. """
177     # Calls the main loss function which handles MPI reduction
178     loss, gU, *_ = loss_and_grad_serial_with_biases(U, W, bu, bi, *rest)
179     return loss, gU # 2-tuple exactly as Catalyst expects
180 def OrthRetraction(U: np.ndarray, V: np.ndarray) -> np.ndarray:
181     """
182     QR-based retraction to the Stiefel / Grassmann manifold.
183     Uses *reduced* QR so it works on NumPy ≥1.26 and CuPy.
184     """
185     # Handle potential zero V vector to avoid QR issues
186     if np.linalg.norm(V) < 1e-12:
187         return U.astype(np.float32)
188
189     # --- FIX: Check for non-finite input ---
190     UV = U + V
191     if not np.isfinite(UV).all():
192         logger.warning("OrthRetraction: Input U+V contains non-finite values. Returning original U.")
193         return U.astype(np.float32)
194     # -----
195
196     try:
197         # --- FIX: Use mode='reduced' ---
198         Q, R_qr = np.linalg.qr(UV, mode='reduced')
199         # -----
200
201         # Ensure Q has the same shape as U
202         if Q.shape[1] < U.shape[1]:
203             pad_width = U.shape[1] - Q.shape[1]
204             Q = np.pad(Q, ((0, 0), (0, pad_width)), mode='constant')
205             logger.warning(f"OrthRetraction: Padded Q due to rank collapse (V norm: {np.linalg.norm(V):.2e})")
206             # Optional: Fix sign ambiguity by matching diagonal of R_qr to be positive
207             # sign_diag = np.sign(np.diag(R_qr))
208             # sign_diag[sign_diag == 0] = 1 # Avoid multiplying by zero
209             # Q = Q @ np.diag(sign_diag)
210             return Q.astype(np.float32)
211         except np.linalg.LinAlgError:
212             logger.warning(f"OrthRetraction: QR decomposition failed (V norm: {np.linalg.norm(V):.2e}). Returning original U.")
213             return U.astype(np.float32)
214         except ValueError as e: # Catch potential value errors from qr
215             logger.error(f"OrthRetraction: ValueError during QR: {e}. Returning original U.")
216             return U.astype(np.float32)
217         except Exception as e: # Catch any other unexpected errors
218             logger.error(f"OrthRetraction failed with unexpected error: {e}")
219             return U.astype(np.float32)
220 # Initialize default values
221 N_users_active, M_movies_active = 0, 0
222 R_train_coo = sparse.coo_matrix((0, 0), dtype=np.float64)
223 R_train_coo_orig = sparse.coo_matrix((0, 0), dtype=np.float64) # For original ratings
224 R_train_csr_orig = sparse.csr_matrix((0,0), dtype=np.float64) # For SoftImpute _matvec
225 R_train_csc_orig = sparse.csc_matrix((0,0), dtype=np.float64) # For SoftImpute _rmatvec
226 ratings_train_orig = np.array([], dtype=np.float64) # Keep original ratings for viz
227 ratings_train_centered = np.array([], dtype=np.float64)
228 mapped_user_ids_train, mapped_movie_ids_train = np.array([], dtype=np.int32), np.array([], dtype=np.int32)
229 user_ids_val_final, movie_ids_val_final, ratings_val_true = (np.array([], dtype=np.int32), np.array([], dtype=np.int32), np.array([], c
230 global_mean_rating = 0.0
231 user_map_global_to_local = {}

```

```

232 movie_map_global_to_local = {}
233 unique_users_train = np.array([], dtype=np.int32)
234 unique_movies_train = np.array([], dtype=np.int32)
235 DATA_AVAILABLE = False
236 user_data_arrays = {} # Precompute user data for ALS/SVRG
237 sampling_prob = None # Initialize sampling probability
238 RSVRG_EPOCH_LEN = 1 # Default epoch length
239
240 ratings_file_path = DATA_DIR / RATINGS_FILENAME
241
242 if DRIVE_MOUNTED and ratings_file_path.is_file():
243     logger.info(f"Loading MovieLens 1M data from: {ratings_file_path}")
244     try:
245         ratings_df = pd.read_csv(
246             ratings_file_path, sep='::', header=None,
247             names=['user_id', 'movie_id', 'rating', 'timestamp'],
248             engine='python', encoding='latin-1'
249         )
250         logger.info(f"Loaded {len(ratings_df)} ratings.")
251         DATA_AVAILABLE = True
252
253         if RATING_LIMIT is not None and RATING_LIMIT > 0 and len(ratings_df) > RATING_LIMIT:
254             logger.info(f"Subsampling ratings from {len(ratings_df)} to {RATING_LIMIT}")
255             ratings_df = ratings_df.sample(n=RATING_LIMIT, random_state=SEED)
256
257         stratify_arg = ratings_df['user_id'] if RATING_LIMIT is None else None
258         if stratify_arg is None and RATING_LIMIT is not None:
259             logger.warning("Stratify is disabled due to RATING_LIMIT being set.")
260         train_df, val_df = train_test_split(
261             ratings_df, test_size=VALIDATION_FRACTION, random_state=SEED, stratify=stratify_arg)
262         logger.info(f"Train size: {len(train_df)}, Validation size: {len(val_df)}")
263
264         user_ids_train_orig = train_df['user_id'].values; movie_ids_train_orig = train_df['movie_id'].values
265         ratings_train_orig = train_df['rating'].values.astype(np.float64)
266         user_ids_val_orig = val_df['user_id'].values; movie_ids_val_orig = val_df['movie_id'].values
267         ratings_val_true = val_df['rating'].values.astype(np.float64) # Keep original for validation
268
269         global_mean_rating = ratings_train_orig.mean()
270         logger.info(f"Global mean rating (training): {global_mean_rating:.4f}")
271
272         unique_users_train, mapped_user_ids_train = np.unique(user_ids_train_orig, return_inverse=True)
273         unique_movies_train, mapped_movie_ids_train = np.unique(movie_ids_train_orig, return_inverse=True)
274         N_users_active = len(unique_users_train); M_movies_active = len(unique_movies_train)
275         user_map_global_to_local = {orig_id: local_id for local_id, orig_id in enumerate(unique_users_train)}
276         movie_map_global_to_local = {orig_id: local_id for local_id, orig_id in enumerate(unique_movies_train)}
277         logger.info(f"Active users in training: {N_users_active}, Active movies in training: {M_movies_active}")
278
279         ratings_train_centered = ratings_train_orig - global_mean_rating
280
281         val_user_mask = np.isin(user_ids_val_orig, unique_users_train)
282         val_movie_mask = np.isin(movie_ids_val_orig, unique_movies_train)
283         val_valid_mask = val_user_mask & val_movie_mask
284         user_ids_val_filt = user_ids_val_orig[val_valid_mask]; movie_ids_val_filt = movie_ids_val_orig[val_valid_mask]
285         ratings_val_true = ratings_val_true[val_valid_mask] # Filter true ratings accordingly
286         user_ids_val_final = np.array([user_map_global_to_local.get(uid, -1) for uid in user_ids_val_filt], dtype=np.int32)
287         movie_ids_val_final = np.array([movie_map_global_to_local.get(mid, -1) for mid in movie_ids_val_filt], dtype=np.int32)
288         valid_map_mask = (user_ids_val_final != -1) & (movie_ids_val_final != -1) # Filter out any potential misses
289         user_ids_val_final = user_ids_val_final[valid_map_mask]; movie_ids_val_final = movie_ids_val_final[valid_map_mask]
290         ratings_val_true = ratings_val_true[valid_map_mask] # Filter again after mapping
291         logger.info(f"Validation pairs mapped to training users/movies: {len(user_ids_val_final)}")
292
293         if ratings_train_centered.size > 0:
294             R_train_coo = sparse.coo_matrix((ratings_train_centered, (mapped_movie_ids_train, mapped_user_ids_train)), shape=(M_movies_
295             R_train_coo.eliminate_zeros()
296             logger.info(f"Built sparse training matrix (Centered) R_train_coo: shape={R_train_coo.shape}, nnz={R_train_coo.nnz}")
297             R_train_coo_orig = sparse.coo_matrix((ratings_train_orig, (mapped_movie_ids_train, mapped_user_ids_train)), shape=(M_movies
298             R_train_coo_orig.eliminate_zeros()
299             R_train_csr_orig = R_train_coo_orig.tocsr(); R_train_csc_orig = R_train_coo_orig.tocsc()
300             logger.info(f"Built sparse training matrix (Original) R_train_coo_orig: shape={R_train_coo_orig.shape}, nnz={R_train_coo_or
301
302         # Precompute user data structures for ALS/SVRG
303         logger.info("Precomputing user data structures...")
304         t_precomp_start = time.time()
305         user_data_arrays = {}
306         for r, c, v in zip(R_train_coo_orig.row, R_train_coo_orig.col, R_train_coo_orig.data):
307             user_data_arrays.setdefault(c, []).append((r, v))
308         for u, rating_list in user_data_arrays.items():

```

```

309         if rating_list:
310             movie_indices_list, rs_list = zip(*rating_list)
311             user_data_arrays[u] = {'movies': np.array(list(movie_indices_list), dtype=np.int32),
312                                     'rs': np.array(list(rs_list), dtype=np.float64)} # Store original ratings
313         logger.info(f"User data precomputation done in {time.time() - t_precomp_start:.2f}s")
314         # Calculate importance sampling weights (consistent across ranks)
315         all_user_indices_global = np.array(list(user_data_arrays.keys()), dtype=np.int32)
316         num_active_users_global = len(all_user_indices_global)
317         user_weights = None; use_importance_sampling = False
318         if num_active_users_global > 0:
319             if RANK_MPI == 0: print("Calculating importance sampling weights...")
320             user_ratings_count = [len(user_data_arrays[u_idx]['movies']) if u_idx in user_data_arrays and 'movies' in user_data_arr
321                                 for u_idx in all_user_indices_global]
322             user_weights_np = np.array(user_ratings_count, dtype=np.float64)
323             sum_weights = user_weights_np.sum()
324             if sum_weights > 1e-9:
325                 user_weights_np /= sum_weights
326                 user_weights = user_weights_np # Probabilities aligned with all_user_indices_global
327                 use_importance_sampling = True
328                 if RANK_MPI == 0: print(f"Importance sampling enabled (weights based on {sum_weights:.0f} ratings).")
329             else:
330                 if RANK_MPI == 0: print("Warning: Cannot compute importance sampling weights. Using uniform.")
331         else:
332             if RANK_MPI == 0: print("No active users, cannot use importance sampling.")
333             sampling_prob = user_weights if use_importance_sampling else None
334             RSVRG_EPOCH_LEN = math.ceil(num_active_users_global / RSVRG_BATCH_SIZE) if num_active_users_global > 0 else 1
335             if RANK_MPI == 0: print(f"RSVRG Epoch Length set to {RSVRG_EPOCH_LEN} batches.")
336         else: logger.error("No training ratings available.")
337
338     except FileNotFoundError: logger.error(f"MovieLens file not found: {ratings_file_path}"); DATA_AVAILABLE = False
339     except Exception as e: logger.error(f"Error processing MovieLens: {e}", exc_info=True); DATA_AVAILABLE = False
340 elif not DRIVE_MOUNTED: logger.error("Google Drive not mounted.")
341 else: logger.error(f"Data directory {DATA_DIR} or ratings file {RATINGS_FILENAME} not found.")
342
343 gc.collect()
344 logger.info("Cell 2: Data Loading and Preprocessing Complete.")
345 logger.info(f"Active Dimensions: M_movies={M_movies_active}, N_users={N_users_active}")
346 logger.info(f"Training Ratings: {R_train_coo.nnz}")
347 logger.info(f"Validation Ratings (for RMSE): {ratings_val_true.size}")
348
349
350 # ===== #
351 # CELL 2.5: Data Visualization
352 # ===== #
353 logger.info("+++ Cell 2.5: Visualizing Loaded Data +++")
354
355 if RANK_MPI == 0: # Only rank 0 should plot
356     if DATA_AVAILABLE and ratings_train_orig.size > 0:
357         plt.style.use('seaborn-v0_8-whitegrid') # Use a nice style
358
359         # 1. Rating Distribution
360         plt.figure(figsize=(10, 4))
361         counts, bins, patches = plt.hist(ratings_train_orig, bins=[0.5, 1.5, 2.5, 3.5, 4.5, 5.5], rwidth=0.8, align='mid', color='skybl
362         bin_centers = 0.5 * (bins[:-1] + bins[1:])
363         for count, x in zip(counts, bin_centers):
364             if count > 0: plt.text(x, count, str(int(count)), ha='center', va='bottom')
365         plt.title('Distribution of Training Ratings (MovieLens 1M Subset)')
366         plt.xlabel('Rating'); plt.ylabel('Frequency')
367         plt.xticks([1, 2, 3, 4, 5]); plt.grid(axis='y', alpha=0.75)
368         plt.tight_layout(); plt.show()
369
370         # 2. Ratings per User
371         user_rating_counts = np.binnedcount(mapped_user_ids_train)
372         plt.figure(figsize=(10, 4))
373         plt.hist(user_rating_counts[user_rating_counts > 0], bins=50, log=True, color='lightcoral', edgecolor='black')
374         plt.title('Distribution of Ratings per User (Training Set)')
375         plt.xlabel('Number of Ratings Given'); plt.ylabel('Number of Users (log scale)')
376         plt.grid(axis='y', alpha=0.75); plt.tight_layout(); plt.show()
377
378         # 3. Ratings per Movie
379         movie_rating_counts = np.binnedcount(mapped_movie_ids_train)
380         plt.figure(figsize=(10, 4))
381         plt.hist(movie_rating_counts[movie_rating_counts > 0], bins=50, log=True, color='lightgreen', edgecolor='black')
382         plt.title('Distribution of Ratings per Movie (Training Set)')
383         plt.xlabel('Number of Ratings Received'); plt.ylabel('Number of Movies (log scale)')
384         plt.grid(axis='y', alpha=0.75); plt.tight_layout(); plt.show()
385         logger.info("Cell 2.5: Data Visualization Complete.")
386     else:

```

```

387         logger.warning("Skipping data visualization as no data was loaded.")
388
389 # ===== #
390 # CELL 3: Model Helpers (CONSOLIDATED)
391 # ===== #
392 logger.info("+++ Cell 3: Defining ALL Model Helpers +++")
393
394 # --- Retraction Factory ---
395 class RetractionFactory:
396     _registry = {}
397     @classmethod
398     def register(cls, name):
399         def decorator(fn): cls._registry[name] = fn; return fn
400         return decorator
401     @classmethod
402     def get(cls, name):
403         if name not in cls._registry: raise KeyError(f"Unknown retraction '{name}'. Available: {list(cls._registry.keys())}")
404         return cls._registry[name]
405 # --- Register Retractions ---
406 @RetractionFactory.register("orthonormal")
407 def _retract_qr(U: np.ndarray, V: np.ndarray) -> np.ndarray:
408     """QR-based retraction."""
409     if np.linalg.norm(V) < 1e-12: return U.astype(np.float32)
410     UV = U + V
411     if not np.isfinite(UV).all(): logger.warning("OrthRetraction: Input U+V non-finite."); return U.astype(np.float32)
412     try:
413         Q, R_qr = np.linalg.qr(UV, mode='reduced') # Use 'reduced'
414         if Q.shape[1] < U.shape[1]:
415             pad_width = U.shape[1] - Q.shape[1]; Q = np.pad(Q, ((0, 0), (0, pad_width)), mode='constant')
416             logger.warning(f"OrthRetraction: Padded Q")
417         return Q.astype(np.float32)
418     except Exception as e: logger.error(f"OrthRetraction failed: {e}"); return U.astype(np.float32)
419 @RetractionFactory.register("cayley")
420 def _retract_cayley(U: np.ndarray, V: np.ndarray, alpha: float = 0.1) -> np.ndarray:
421     """ Simple Cayley approx using QR of ambient step. """
422     return _retract_qr(U, alpha * V)
423 @RetractionFactory.register("projection")
424 def _retract_projection(U: np.ndarray, V: np.ndarray) -> np.ndarray:
425     """ Projection (polar decomposition) retraction. """
426     U64 = U.astype(np.float64, copy=False); V64 = V.astype(np.float64, copy=False)
427     Z = U64 + V64; G = Z.T @ Z
428     try:
429         s, P = np.linalg.eigh(G); s_safe = np.maximum(s, 1e-12)
430         s_inv_sqrt = 1.0 / np.sqrt(s_safe); G_mhalf = P @ np.diag(s_inv_sqrt) @ P.T
431         result = (Z @ G_mhalf).astype(np.float32)
432         if result.shape != U.shape: logger.warning(f"Projection Retraction Warning: Shape mismatch. Falling back to QR."); return _retract_qr(U, V)
433         return result
434     except Exception as e: logger.warning(f"Projection Retraction Warning: {e}. Falling back to QR."); return _retract_qr(U, V)
435 # --- Get the chosen retraction function ---
436 R_fn = RetractionFactory.get(RETRACTION_NAME)
437 if RANK_MPI == 0: logger.info(f"Using Retraction: {RETRACTION_NAME}")
438
439 # --- Other Manifold Helpers ---
440 def ProjTangent(U: np.ndarray, G: np.ndarray) -> np.ndarray:
441     """Project G onto tangent space at U (Grassmann)."""
442     return (G - U @ (U.T @ G)).astype(np.float32)
443 def LogMapApprox(U_base: np.ndarray, U_target: np.ndarray) -> np.ndarray:
444     """Approximate inverse retraction (log map)."""
445     return ProjTangent(U_base, U_target - U_base)
446 def RegularizeGradChordalApprox(U: np.ndarray, U_old: np.ndarray, kappa: float) -> np.ndarray:
447     """Approximate gradient of distance regularization term."""
448     U = U.astype(np.float32); U_old = U_old.astype(np.float32);
449     if REG_DISTANCE == "euclid": S = U.T @ U_old; grad_ambient = U @ (S - S.T); return kappa * ProjTangent(U, grad_ambient)
450     elif REG_DISTANCE == "retraction": v = LogMapApprox(U, U_old); return -kappa * v
451     else: raise ValueError(f"Unknown REG_DISTANCE type: {REG_DISTANCE}")
452
453 # --- RMSE Evaluation ---
454 def evaluate_rmse_with_biases(
455     U: np.ndarray, W: np.ndarray,
456     user_bias: np.ndarray, movie_bias: np.ndarray, global_mean: float,
457     probe_users_mapped: np.ndarray, probe_movies_mapped: np.ndarray, probe_ratings_true: np.ndarray # Now contains true ratings
458 ) -> float:
459     """Computes RMSE on the validation set including bias terms and clamping."""
460     if probe_ratings_true.size == 0: return np.nan # Check if validation set is empty
461     U = U.astype(np.float64, copy=False); W = W.astype(np.float64, copy=False)
462     user_bias = user_bias.astype(np.float64, copy=False); movie_bias = movie_bias.astype(np.float64, copy=False)
463     local_sum_sq_err = 0.0; local_count = 0

```

```

464     try:
465         if M_movies_active == 0 or N_users_active == 0: return np.nan
466         if probe_movies_mapped.size > 0 and (probe_movies_mapped.max() >= M_movies_active or probe_movies_mapped.min() < 0): return np.
467         if probe_users_mapped.size > 0 and (probe_users_mapped.max() >= N_users_active or probe_users_mapped.min() < 0): return np.nan
468         dot_prods = np.array([np.dot(U[m, :], W[:, u]) for m, u in zip(probe_movies_mapped, probe_users_mapped)], dtype=np.float64)
469         preds_raw = global_mean + user_bias[probe_users_mapped] + movie_bias[probe_movies_mapped] + dot_prods
470         preds_clamped = np.clip(preds_raw, 1.0, 5.0)
471         if not np.isfinite(preds_clamped).all(): preds_clamped = np.nan_to_num(preds_clamped, nan=global_mean)
472         if not np.isfinite(probe_ratings_true).all(): probe_ratings_true = np.nan_to_num(probe_ratings_true)
473         squared_errors = (preds_clamped - probe_ratings_true)**2
474         local_sum_sq_err = np.sum(squared_errors)
475         local_count = len(squared_errors)
476     except IndexError as e: logger.error(f"IndexError during biased RMSE: {e}"); return np.nan
477     except Exception as e: logger.error(f"Error during biased RMSE: {e}"); return np.nan
478     # --- MPI Reduction for RMSE ---
479     if COMM and SIZE_MPI > 1:
480         global_sum_sq_err_buf = np.array(local_sum_sq_err, dtype=np.float64); global_count_buf = np.array(local_count, dtype=np.int64)
481         global_sum_sq_err = np.array(0.0, dtype=np.float64); global_count = np.array(0, dtype=np.int64)
482         COMM.Allreduce(global_sum_sq_err_buf, global_sum_sq_err, op=MPI.SUM); COMM.Allreduce(global_count_buf, global_count, op=MPI.SUM)
483         if global_count > 0: mean_squared_error = global_sum_sq_err / global_count
484         else: return np.nan
485     else: # Serial case
486         if local_count > 0: mean_squared_error = local_sum_sq_err / local_count
487         else: return np.nan
488     mean_squared_error = max(0.0, mean_squared_error); rmse = np.sqrt(mean_squared_error)
489     return float(rmse) if np.isfinite(rmse) else np.nan
490
491 # --- RMSE Helper for SoftImpute (No Biases) ---
492 def evaluate_rmse_low_rank(U, S, V, probe_movies_mapped, probe_users_mapped, probe_ratings_true, use_gpu=False):
493     """Computes RMSE for low-rank model X = USV^T against true ratings."""
494     if probe_ratings_true.size == 0: return np.nan
495     xp = cp if use_gpu else np
496     try:
497         if M_movies_active == 0 or N_users_active == 0: return np.nan
498         if probe_movies_mapped.size > 0 and (probe_movies_mapped.max() >= M_movies_active or probe_movies_mapped.min() < 0): return np.
499         if probe_users_mapped.size > 0 and (probe_users_mapped.max() >= N_users_active or probe_users_mapped.min() < 0): return np.nan
500         U_dev = xp.asarray(U); S_dev = xp.asarray(S); V_dev = xp.asarray(V)
501         probe_movies_dev = xp.asarray(probe_movies_mapped); probe_users_dev = xp.asarray(probe_users_mapped)
502         probe_ratings_dev = xp.asarray(probe_ratings_true)
503         term2 = S_dev * V_dev[probe_users_dev, :]
504         preds_raw = xp.sum(U_dev[probe_movies_dev, :] * term2, axis=1)
505         preds_clamped = xp.clip(preds_raw, 1.0, 5.0)
506         if not xp.isfinite(preds_clamped).all(): preds_clamped = xp.nan_to_num(preds_clamped, nan=3.0)
507         if not xp.isfinite(probe_ratings_dev).all(): probe_ratings_dev = xp.nan_to_num(probe_ratings_dev)
508         mse_dev = xp.mean((preds_clamped - probe_ratings_dev)**2)
509         mse = float(cp.asnumpy(mse_dev) if use_gpu else mse_dev)
510         rmse = np.sqrt(mse) if mse >= 0 else np.nan
511     except IndexError as e: logger.error(f"IndexError during low-rank RMSE: {e}"); return np.nan
512     except Exception as e: logger.error(f"Error during low-rank RMSE: {e}"); return np.nan
513     return float(rmse) if np.isfinite(rmse) else np.nan
514
515 # --- Initialization ---
516 def initialize_factors_and_biases(M: int, N: int, R: int, rng: Generator, scale: float) -> Tuple[np.ndarray, np.ndarray, np.ndarray, np
517     """Initializes U, W, user_bias, movie_bias."""
518     U = None; W = None; user_bias = None; movie_bias = None
519     if RANK_MPI == 0:
520         U = rng.standard_normal(size=(M, R)).astype(np.float64) * scale
521         W = rng.standard_normal(size=(R, N)).astype(np.float64) * scale
522         user_bias = np.zeros(N, dtype=np.float64)
523         movie_bias = np.zeros(M, dtype=np.float64)
524         if M >= R: U_orth, _ = np.linalg.qr(U, mode='reduced'); U = U_orth.astype(np.float64)
525         else: logger.warning(f"M ({M}) < R ({R}). Cannot orthonormalize U.")
526     if COMM and SIZE_MPI > 1:
527         if RANK_MPI != 0: U = np.empty((M, R), dtype=np.float64); W = np.empty((R, N), dtype=np.float64); user_bias = np.empty(N, dtype
528         COMM.Bcast(U, root=0); COMM.Bcast(W, root=0); COMM.Bcast(user_bias, root=0); COMM.Bcast(movie_bias, root=0)
529     return U, W, user_bias, movie_bias
530
531 # --- Initial State Recorder ---
532 def record_initial_state_biased(U, W, user_bias, movie_bias, loss_args_biased, eval_args_biased):
533     """Computes and logs initial state for biased models."""
534     current_loss, gU0, gW0, gBu0, gBi0 = loss_and_grad_serial_with_biases(U, W, user_bias, movie_bias, *loss_args_biased)
535     current_rmse = evaluate_rmse_with_biases(U, W, user_bias, movie_bias, *eval_args_biased)
536     gU_proj_0 = ProjTangent(U, gU0)
537     grad_norm_U_riemann = np.linalg.norm(gU_proj_0)
538     grad_norm_W = np.linalg.norm(gW0); grad_norm_Bu = np.linalg.norm(gBu0); grad_norm_Bi = np.linalg.norm(gBi0)
539     if RANK_MPI == 0: logger.info(
540         f"Epoch 00 (Init): Loss={current_loss:.4e}, RMSE={current_rmse:.4f}, "
541         f"||grad||_U={grad_norm_U_riemann:.2e}, ||grad||_W={grad_norm_W:.2e}, "

```

```

542         f"||gBu||={grad_norm_Bu:.2e}, ||gBi||={grad_norm_Bi:.2e}"
543     )
544     if not np.isfinite(current_loss): raise ValueError("Initial loss is not finite.")
545     return current_loss, current_rmse, gU0, gW0, gBu0, gBi0
546
547 # --- Armijo Line Search ---
548 def ArmijoLineSearchRiemannian(
549     U: np.ndarray, G_euclidean: np.ndarray, loss_args: tuple, current_loss: float,
550     lr_init: float, beta: float, sigma: float, max_ls_iter: int = 20
551 ) -> Tuple[float, np.ndarray, float]:
552     """Performs Armijo line search using retraction."""
553     lr = lr_init
554     G_proj = ProjTangent(U, G_euclidean)
555     G_proj_norm_sq = np.linalg.norm(G_proj)**2
556     if G_proj_norm_sq < 1e-14: return 0.0, U, current_loss
557     for ls_iter in range(max_ls_iter):
558         step_vec = -lr * G_proj
559         U_next = R_fn(U, step_vec) # Use chosen retraction
560         if not np.isfinite(U_next).all(): lr *= beta; continue
561         try:
562             W_ls, ub_ls, mb_ls, *rest_args = loss_args
563             loss_next, _, _, _ = loss_and_grad_serial_with_biases(U_next, W_ls, ub_ls, mb_ls, *rest_args)
564             except Exception as e: logger.error(f"Armijo LS Error: {e}"); return 0.0, U, current_loss
565             if not np.isfinite(loss_next): lr *= beta; continue
566             required_decrease = sigma * lr * G_proj_norm_sq
567             actual_decrease = current_loss - loss_next
568             if actual_decrease >= required_decrease - 1e-9: return lr, U_next, loss_next
569             lr *= beta
570             if lr < 1e-14: break
571     logger.debug("Armijo LS failed."); return 0.0, U, current_loss
572
573 # --- Adaptive Stopping Check ---
574 def should_stop_subproblem(G_proj, step_vec):
575     """Return True if both criteria are already small."""
576     grad_norm_proj = np.linalg.norm(G_proj)
577     step_norm = np.linalg.norm(step_vec)
578     stop = (grad_norm_proj < ETA_GRAD and step_norm < ETA_DIST)
579     return stop
580
581 # --- Adaptive Kappa Update ---
582 def update_kappa_adaptive(kappa_prev, h_hist, dist_hist, U_local,
583                           gamma=2.0, window=3,
584                           kappa_min=1e-4, kappa_max=1e12):
585     """ Adaptive kappa update using local curvature estimate. """
586     if U_local.shape[1] == 0: return kappa_min # Handle empty matrix case
587     v = GLOBAL_RNG.standard_normal(size=(U_local.shape[1], 1)).astype(U_local.dtype)
588     v /= np.linalg.norm(v) + 1e-12
589     U_local_64 = U_local.astype(np.float64); v_64 = v.astype(np.float64)
590     lambda_max_sq = 0.0
591     for _ in range(2): # 2 power iterations on U^T U
592         Av = U_local_64.T @ (U_local_64 @ v_64)
593         lambda_max_sq = v_64.T @ Av
594         v_norm = np.linalg.norm(Av); v_64 = Av / (v_norm + 1e-12)
595     L_local = np.sqrt(max(0.0, lambda_max_sq.item()))
596     target_ratio = 0.9; target = target_ratio * L_local
597     kappa_new = np.clip(target, kappa_min, kappa_max)
598     return float(kappa_new)
599
600 # --- OT Demo Helper ---
601 def run_barycentre_demo(n_grid=200, reg=1e-1, rng_seed=0):
602     """ POT demo: 3 one-dimensional Gaussians -> entropic Wasserstein barycenter """
603     if not OT_AVAILABLE: return None
604     grid = np.linspace(-8.0, 8.0, n_grid)
605     M = ot.dist(grid.reshape(-1, 1), grid.reshape(-1, 1)) ** 2
606     means = np.array([-3.0, 0.0, 3.0]); sigmas = np.array([0.5, 1.0, 0.7])
607     sources = np.vstack([np.exp(-0.5 * ((grid - m) / s) ** 2) / (s * np.sqrt(2 * np.pi)) for m, s in zip(means, sigmas)]).T
608     sources /= sources.sum(axis=0, keepdims=True)
609     bary, log = ot.bregman.barycenter(sources, M, reg, weights=None, numItermax=1000, stopThr=1e-7, log=True)
610     return {'grid': grid, 'sources': sources, 'barycenter': bary, 'log': log}
611
612
613 logger.info("Cell 3: Model Helpers Defined.")
614
615 # =====
616 # CELL 4: Non-Convex Solvers (SVRG, ALS, Euclidean GD) - Renumbered
617 # =====
618 logger.info("+++ Cell 4: Defining Non-Convex Solvers +++")

```



```

619 # --- Loss/Gradient Functions ---
620 def loss_and_grad_serial_with_biases(
621     U: np.ndarray, W: np.ndarray, user_bias: np.ndarray, movie_bias: np.ndarray,
622     global_mean: float,
623     rows_idx: np.ndarray, cols_idx: np.ndarray, vals_true_centered: np.ndarray, # Centered ratings
624     n_movies_func: int, n_users_func: int, rank_func: int,
625     lambda_sq_func: float, lambda_bias_func: float
626 ) -> Tuple[float, np.ndarray, np.ndarray, np.ndarray, np.ndarray]:
627     """ Computes loss and gradients for U, W, user_bias, movie_bias. """
628     # ... (implementation from v11) ...
629     U = U.astype(np.float64, copy=False); W = W.astype(np.float64, copy=False)
630     user_bias = user_bias.astype(np.float64, copy=False); movie_bias = movie_bias.astype(np.float64, copy=False)
631     if vals_true_centered.size == 0: return 0.0, np.zeros_like(U), np.zeros_like(W), np.zeros_like(user_bias), np.zeros_like(movie_bias)
632     try:
633         W_cols = W[:, cols_idx]; U_rows = U[rows_idx, :]
634         dot_prods = np.sum(U_rows * W_cols.T, axis=1)
635         preds_residual = user_bias[cols_idx] + movie_bias[rows_idx] + dot_prods
636     except IndexError as e: logger.error(f"Indexing error in loss_and_grad_serial_with_biases - {e}"); raise
637     valid_mask = np.isfinite(preds_residual) & np.isfinite(vals_true_centered)
638     if not np.all(valid_mask):
639         logger.warning(f"Filtering {np.sum(~valid_mask)} non-finite values in loss_and_grad_serial_with_biases.")
640         rows_idx_filt = rows_idx[valid_mask]; cols_idx_filt = cols_idx[valid_mask]
641         vals_true_filt = vals_true_centered[valid_mask]; preds_filt = preds_residual[valid_mask]
642         if preds_filt.size == 0: return np.inf, np.zeros_like(U), np.zeros_like(W), np.zeros_like(user_bias), np.zeros_like(movie_bias)
643     else:
644         rows_idx_filt, cols_idx_filt, vals_true_filt, preds_filt = rows_idx, cols_idx, vals_true_centered, preds_residual
645     errors = preds_filt - vals_true_filt
646     loss_obs = 0.5 * np.sum(errors**2)
647     loss_reg_U = 0.5 * lambda_sq_func * np.sum(U**2); loss_reg_W = 0.5 * lambda_sq_func * np.sum(W**2)
648     loss_reg_bu = 0.5 * lambda_bias_func * np.sum(user_bias**2); loss_reg_bi = 0.5 * lambda_bias_func * np.sum(movie_bias**2)
649     total_loss = loss_obs + loss_reg_U + loss_reg_W + loss_reg_bu + loss_reg_bi
650     E_sparse = sparse.csr_matrix((errors, (rows_idx_filt, cols_idx_filt)), shape=(n_movies_func, n_users_func))
651     E_sparse_csc = E_sparse.tocsc()
652     grad_U = E_sparse @ W.T + lambda_sq_func * U
653     grad_W = U.T @ E_sparse_csc + lambda_sq_func * W
654     grad_user_bias = np.array(E_sparse.sum(axis=0)).flatten() + lambda_bias_func * user_bias
655     grad_movie_bias = np.array(E_sparse.sum(axis=1)).flatten() + lambda_bias_func * movie_bias
656     if not np.isfinite(grad_U).all(): grad_U = np.nan_to_num(grad_U)
657     if not np.isfinite(grad_W).all(): grad_W = np.nan_to_num(grad_W)
658     if not np.isfinite(grad_user_bias).all(): grad_user_bias = np.nan_to_num(grad_user_bias)
659     if not np.isfinite(grad_movie_bias).all(): grad_movie_bias = np.nan_to_num(grad_movie_bias)
660     if not np.isfinite(total_loss): total_loss = np.inf
661     return float(total_loss), grad_U.astype(np.float64), grad_W.astype(np.float64), grad_user_bias.astype(np.float64), grad_movie_bias.
662
663 def gradient_batch_with_biases(
664     U: np.ndarray, W: np.ndarray, user_bias: np.ndarray, movie_bias: np.ndarray,
665     indices: np.ndarray, # Indices into GLOBAL triplets
666     rows_idx: np.ndarray, cols_idx: np.ndarray, vals_true_centered: np.ndarray, # Centered ratings
667     n_ratings_total: int,
668     lambda_sq_func: float, lambda_bias_func: float
669 ) -> Tuple[np.ndarray, np.ndarray, np.ndarray, np.ndarray]:
670     """ Computes average Euclidean gradient over a BATCH of ratings, including biases. """
671     U = U.astype(np.float64, copy=False)
672     W = W.astype(np.float64, copy=False)
673     user_bias = user_bias.astype(np.float64, copy=False)
674     movie_bias = movie_bias.astype(np.float64, copy=False)
675     batch_size = len(indices)
676     if batch_size == 0:
677         return np.zeros_like(U), np.zeros_like(W), np.zeros_like(user_bias), np.zeros_like(movie_bias)
678
679     # Get data for the batch
680     batch_rows = rows_idx[indices]
681     batch_cols = cols_idx[indices]
682     batch_vals_centered = vals_true_centered[indices]
683
684     # Get corresponding factors and biases
685     try:
686         U_batch = U[batch_rows, :] # Shape (B, R)
687         W_batch = W[:, batch_cols] # Shape (R, B)
688         user_bias_batch = user_bias[batch_cols] # Shape (B,)
689         movie_bias_batch = movie_bias[batch_rows] # Shape (B,)
690     except IndexError as e:
691         logger.error(f"Indexing error in gradient_batch_with_biases - {e}")
692         raise
693
694     # Predict residual for the batch
695     preds_batch_residual = user_bias_batch + movie_bias_batch + np.sum(U_batch * W_batch.T, axis=1)

```

```

696
697 # Calculate errors for the batch
698 errors_batch = preds_batch_residual - batch_vals_centered # Shape (B,)
699
700 # Calculate gradient terms using sparse matrix approach
701 E_sparse_batch = sparse.csr_matrix((errors_batch, (batch_rows, batch_cols)),
702                                     shape=(U.shape[0], W.shape[1]))
703
704 # Average gradient over the batch
705 grad_U_batch = (E_sparse_batch @ W.T) / batch_size + lambda_sq_func * U
706 grad_W_batch = (U.T @ E_sparse_batch.tocsc()) / batch_size + lambda_sq_func * W
707
708 # Compute bias gradients (need to average errors per user/movie in batch)
709 # This requires accumulating errors per user/movie index present in the batch
710 grad_user_bias_batch = np.zeros_like(user_bias)
711 grad_movie_bias_batch = np.zeros_like(movie_bias)
712 np.add.at(grad_user_bias_batch, batch_cols, errors_batch) # Accumulate errors by user index
713 np.add.at(grad_movie_bias_batch, batch_rows, errors_batch) # Accumulate errors by movie index
714
715 grad_user_bias_batch = grad_user_bias_batch / batch_size + lambda_bias_func * user_bias
716 grad_movie_bias_batch = grad_movie_bias_batch / batch_size + lambda_bias_func * movie_bias
717
718 # Handle potential non-finite values
719 if not np.isfinite(grad_U_batch).all(): grad_U_batch = np.nan_to_num(grad_U_batch)
720 if not np.isfinite(grad_W_batch).all(): grad_W_batch = np.nan_to_num(grad_W_batch)
721 if not np.isfinite(grad_user_bias_batch).all(): grad_user_bias_batch = np.nan_to_num(grad_user_bias_batch)
722 if not np.isfinite(grad_movie_bias_batch).all(): grad_movie_bias_batch = np.nan_to_num(grad_movie_bias_batch)
723
724 return grad_U_batch.astype(np.float64), grad_W_batch.astype(np.float64), grad_user_bias_batch.astype(np.float64), grad_movie_bias_t
725
726 # --- SVRG Solver ---
727 # --- SVRG Solver with Biases ---
728 def run_non_convex_svrg_with_biases(
729     R_train_coo: sparse.coo_matrix, # Contains centered ratings
730     global_mean: float,
731     probe_users_mapped: np.ndarray, # Mapped probe indices
732     probe_movies_mapped: np.ndarray,
733     probe_ratings_true: np.ndarray, # Original probe ratings
734     N_users_active: int,
735     M_movies_active: int,
736     rank_local: int,
737     n_epochs: int,
738     inner_lr: float, # Base inner learning rate
739     batch_size: int,
740     lam_sq: float,
741     lam_bias: float,
742     rng: Generator,
743     init_scale: float = INIT_SCALE_NON_CONVEX,
744     max_grad_norm: float = GRAD_CLIP_THRESHOLD
745 ) -> Dict[str, List]:
746     """
747     Runs SVRG for non-convex UW factorization including bias terms.
748     Uses decaying LR and gradient clipping.
749     """
750     logger.info("Starting Non-Convex SVRG Solver with Biases...")
751     # Initialize factors and biases
752     U, W, user_bias, movie_bias = initialize_factors_and_biases(
753         M_movies_active, N_users_active, rank_local, rng, init_scale
754     )
755
756     hist_loss = []
757     hist_rmse = []
758     hist_time = []
759     hist_gU_norm, hist_gW_norm, hist_gBu_norm, hist_gBi_norm = [], [], [], []
760
761     start_time = time.time()
762
763     # Use mapped indices and centered ratings for training
764     train_rows = R_train_coo.row
765     train_cols = R_train_coo.col
766     train_vals_centered = R_train_coo.data
767     n_ratings_total = R_train_coo.nnz
768
769     if n_ratings_total == 0:
770         logger.error("No training ratings available.")
771         return {'loss': [], 'rmse': [], 'time': [], 'gU_norm': [], 'gW_norm': [], 'gBu_norm': [], 'gBi_norm': [], 'U': None, 'W': None,
772
773     # Initial evaluation

```

```

774     try:
775         loss0, gU0, gW0, gBu0, gBi0 = loss_and_grad_serial_with_biases(
776             U, W, user_bias, movie_bias, global_mean,
777             train_rows, train_cols, train_vals_centered,
778             M_movies_active, N_users_active, rank_local, lam_sq, lam_bias
779         )
780         rmse0 = evaluate_rmse_with_biases(
781             U, W, user_bias, movie_bias, global_mean,
782             probe_users_mapped, probe_movies_mapped, probe_ratings_true
783         )
784         hist_loss.append(loss0)
785         hist_rmse.append(rmse0)
786         hist_time.append(time.time() - start_time)
787         hist_gU_norm.append(np.linalg.norm(gU0))
788         hist_gW_norm.append(np.linalg.norm(gW0))
789         hist_gBu_norm.append(np.linalg.norm(gBu0))
790         hist_gBi_norm.append(np.linalg.norm(gBi0))
791         logger.info(
792             f"Epoch 00 (Init): Loss={loss0:.4e}, RMSE={rmse0:.4f}, "
793             f"||gU||={hist_gU_norm[-1]:.2e}, ||gW||={hist_gW_norm[-1]:.2e}, "
794             f"||gBu||={hist_gBu_norm[-1]:.2e}, ||gBi||={hist_gBi_norm[-1]:.2e}"
795         )
796     except Exception as e:
797         logger.error(f"Error during initial evaluation: {e}", exc_info=True)
798         return {'loss': [], 'rmse': [], 'time': [], 'gU_norm': [], 'gW_norm': [], 'gBu_norm': [], 'gBi_norm': [], 'U': None, 'W': None,
799
800 # Main SVRG Loop
801 for epoch in range(1, n_epochs + 1):
802     epoch_start_time = time.time()
803     logger.info(f"--- Starting Epoch {epoch:02d} ---")
804     # --- Use Exponential Decay for Learning Rate (FIX 4) ---
805     lr_epoch = inner_lr * (0.9**(epoch - 1)) # Exponential decay
806     logger.info(f"Using lr = {lr_epoch:.2e} this epoch")
807     # -----
808
809     # Compute anchor gradient
810     logger.info(f"Epoch {epoch:02d}: Computing anchor gradient...")
811     anchor_start_time = time.time()
812     try:
813         loss_anchor, gU_anchor, gW_anchor, gBu_anchor, gBi_anchor = loss_and_grad_serial_with_biases(
814             U, W, user_bias, movie_bias, global_mean,
815             train_rows, train_cols, train_vals_centered,
816             M_movies_active, N_users_active, rank_local, lam_sq, lam_bias
817         )
818         logger.info(f"Epoch {epoch:02d}: Anchor gradient computed in {time.time() - anchor_start_time:.2f}s.")
819     except Exception as e:
820         logger.error(f"Error computing anchor gradient at epoch {epoch}: {e}")
821         break
822
823     U_epoch_start, W_epoch_start = U.copy(), W.copy()
824     user_bias_epoch_start, movie_bias_epoch_start = user_bias.copy(), movie_bias.copy()
825
826     # Inner loop
827     # --- Use Full Inner Pass (FIX 5) ---
828     num_inner_steps = max(1, (n_ratings_total // batch_size) // SVRG_INNER_STEPS_DIVISOR)
829     logger.info(f"Epoch {epoch:02d}: Starting inner loop with {num_inner_steps} steps...")
830     inner_loop_start_time = time.time()
831
832     for inner_step in range(num_inner_steps):
833         batch_indices = rng.choice(n_ratings_total, size=batch_size, replace=False)
834         try:
835             gU_curr, gW_curr, gBu_curr, gBi_curr = gradient_batch_with_biases(
836                 U, W, user_bias, movie_bias, batch_indices,
837                 train_rows, train_cols, train_vals_centered,
838                 n_ratings_total, lam_sq, lam_bias)
839             gU_anch, gW_anch, gBu_anch, gBi_anch = gradient_batch_with_biases(
840                 U_epoch_start, W_epoch_start, user_bias_epoch_start, movie_bias_epoch_start,
841                 batch_indices, train_rows, train_cols, train_vals_centered,
842                 n_ratings_total, lam_sq, lam_bias)
843         except Exception as e:
844             logger.error(f"Error computing stochastic gradient: {e}")
845             continue
846
847         # Variance-reduced gradients
848         gU_vr = gU_curr - gU_anch + gU_anchor
849         gW_vr = gW_curr - gW_anch + gW_anchor
850         gBu_vr = gBu_curr - gBu_anch + gBu_anchor

```

```

851     gBi_vr = gBi_curr - gBi_anch + gBi_anchor
852
853     # Gradient clipping
854     gU_norm = np.linalg.norm(gU_vr); gW_norm = np.linalg.norm(gW_vr)
855     gBu_norm = np.linalg.norm(gBu_vr); gBi_norm = np.linalg.norm(gBi_vr)
856     if gU_norm > max_grad_norm: gU_vr *= (max_grad_norm / gU_norm)
857     if gW_norm > max_grad_norm: gW_vr *= (max_grad_norm / gW_norm)
858     if gBu_norm > max_grad_norm: gBu_vr *= (max_grad_norm / gBu_norm)
859     if gBi_norm > max_grad_norm: gBi_vr *= (max_grad_norm / gBi_norm)
860
861     # Update factors and biases
862     U -= lr_epoch * gU_vr
863     W -= lr_epoch * gW_vr
864     user_bias -= lr_epoch * gBu_vr
865     movie_bias -= lr_epoch * gBi_vr
866
867     if (inner_step + 1) % 5000 == 0: # Log less frequently for full inner pass
868         logger.info(f"Epoch {epoch:02d}: Inner step {inner_step+1}/{num_inner_steps} done.")
869
870     logger.info(f"Epoch {epoch:02d}: Inner loop finished in {time.time() - inner_loop_start_time:.2f}s.")
871
872     # Evaluate after epoch
873     logger.info(f"Epoch {epoch:02d}: Evaluating loss and RMSE...")
874     eval_start_time = time.time()
875     try:
876         loss_k, gU_k, gW_k, gBu_k, gBi_k = loss_and_grad_serial_with_biases(
877             U, W, user_bias, movie_bias, global_mean,
878             train_rows, train_cols, train_vals_centered,
879             M_movies_active, N_users_active, rank_local, lam_sq, lam_bias
880         )
881         if not np.isfinite(loss_k):
882             logger.error(f"Epoch {epoch:02d}: Loss became non-finite ({loss_k}). Stopping.")
883             hist_loss.append(np.nan); hist_rmse.append(np.nan); hist_time.append(time.time() - start_time)
884             hist_gU_norm.append(np.nan); hist_gW_norm.append(np.nan); hist_gBu_norm.append(np.nan); hist_gBi_norm.append(np.nan)
885             break
886
887         rmse_k = evaluate_rmse_with_biases(
888             U, W, user_bias, movie_bias, global_mean,
889             probe_users_mapped, probe_movies_mapped, probe_ratings_true
890         )
891         hist_loss.append(loss_k); hist_rmse.append(rmse_k)
892         hist_time.append(time.time() - start_time)
893         hist_gU_norm.append(np.linalg.norm(gU_k)); hist_gW_norm.append(np.linalg.norm(gW_k))
894         hist_gBu_norm.append(np.linalg.norm(gBu_k)); hist_gBi_norm.append(np.linalg.norm(gBi_k))
895
896         logger.info(f"Epoch {epoch:02d}: Eval done in {time.time() - eval_start_time:.2f}s.")
897         logger.info(
898             f"Loss={loss_k:.4e}, RMSE={rmse_k:.4f}, "
899             f"||gU||={hist_gU_norm[-1]:.2e}, ||gW||={hist_gW_norm[-1]:.2e}, "
900             f"||gBu||={hist_gBu_norm[-1]:.2e}, ||gBi||={hist_gBi_norm[-1]:.2e}"
901         )
902     except Exception as e:
903         logger.error(f"Error during evaluation at epoch {epoch}: {e}", exc_info=True)
904         hist_loss.append(np.nan); hist_rmse.append(np.nan); hist_time.append(time.time() - start_time)
905         hist_gU_norm.append(np.nan); hist_gW_norm.append(np.nan); hist_gBu_norm.append(np.nan); hist_gBi_norm.append(np.nan)
906         break
907
908     logger.info(f"--- Epoch {epoch:02d} finished in {time.time() - epoch_start_time:.2f}s ---")
909
910     logger.info("Non-Convex SVRG Solver with Biases Finished.")
911     return {
912         'loss': hist_loss, 'rmse': hist_rmse, 'time': hist_time,
913         'gU_norm': hist_gU_norm, 'gW_norm': hist_gW_norm,
914         'gBu_norm': hist_gBu_norm, 'gBi_norm': hist_gBi_norm,
915         'U': U, 'W': W, 'bu': user_bias, 'bi': movie_bias
916     }
917
918
919 # --- ALS Solver ---
920
921 def W_closed_efficient(U, N_users, N_movies, user_indices=None):
922     # Solves for W for a subset of users (local computation)
923     U = U.astype(np.float32, copy=False);
924     target_users = user_indices if user_indices is not None else user_data_arrays.keys()
925     W_subset = {} # Use dict if only computing for subset
926     I_r_lam_sq = (LAM_SQ * I_r).astype(np.float32) # lambda^2 * I
927
928     for u in target_users:

```

```

929         if u not in user_data_arrays: continue
930         data = user_data_arrays[u]
931         movie_indices = data['movies']; rs_t = data['rs']
932         if movie_indices.size == 0: continue
933         # Check bounds before indexing U
934         if movie_indices.max() >= U.shape[0] or movie_indices.min() < 0:
935             # if RANK_MPI == 0: print(f"Warning: Invalid movie indices for user {u}. Skipping.")
936             continue
937         U_k = U[movie_indices, :]
938         A = U_k.T @ U_k + I_r_lam_sq
939         B = U_k.T @ rs_t
940         A = A.astype(np.float32); B = B.astype(np.float32)
941         try:
942             w_u = np.linalg.solve(A.astype(np.float64), B.astype(np.float64)).astype(np.float32)
943         except np.linalg.LinAlgError:
944             # if RANK_MPI == 0: print(f"Warning: np.linalg.solve failed for user {u}. Using pseudo-inverse.")
945             try:
946                 w_u = (np.linalg.pinv(A.astype(np.float64)) @ B.astype(np.float64)).astype(np.float32)
947             except np.linalg.LinAlgError:
948                 if RANK_MPI == 0: print(f"ERROR: Pseudo-inverse also failed for user {u}. Returning zero vector.")
949                 w_u = np.zeros(RANK, dtype=np.float32) # Return zero vector if fails completely
950             except Exception as e_pinv:
951                 if RANK_MPI == 0: print(f"ERROR: Unknown error in pseudo-inverse for user {u}: {e_pinv}. Returning zero vector.")
952                 w_u = np.zeros(RANK, dtype=np.float32)
953
954         if user_indices is not None:
955             W_subset[u] = w_u
956         else:
957             if 'W' not in locals(): W = np.zeros((RANK, N_users), dtype=np.float32)
958             if 0 <= u < W.shape[1]: # Check user index bound for W
959                 W[:, u] = w_u
960             # else: # This shouldn't happen if N_users is correct
961             #     if RANK_MPI == 0: print(f"Warning: User index {u} out of bounds for W (shape {W.shape}).")
962
963     if user_indices is not None:
964         return W_subset # Return dict
965     else:
966         if 'W' not in locals():
967             # if RANK_MPI == 0: print("Warning: W_closed_efficient called with no active users? Returning empty W.")
968             return np.zeros((RANK, N_users), dtype=np.float32)
969         # W should be filled now
970         if not np.isfinite(W).all():
971             if RANK_MPI == 0: print("Warning: Non-finite values found in computed W matrix. Clamping.")
972             W = np.nan_to_num(W, nan=0.0, posinf=0.0, neginf=0.0) # Clamp non-finite to zero
973         assert W.shape == (RANK, N_users);
974         return W # Return full W matrix
975
976
977
978 def update_user_factors(
979     R_train_coo_csc: sparse.csc_matrix, # Centered ratings, CSC format
980     U: np.ndarray,
981     user_bias: np.ndarray,
982     movie_bias: np.ndarray,
983     lambda_sq: float,
984     rank: int,
985     N_users: int
986 ) -> np.ndarray:
987     """Solves for W (user factors) fixing U and biases."""
988     M = U.shape[0]
989     W = np.zeros((rank, N_users), dtype=np.float64)
990     # Precompute U^T U + lambda*I (used in the denominator)
991     # Note: This is used inside the loop per user based on specific movies U_j
992     # UtU = U.T @ U + lambda_sq * np.eye(rank, dtype=np.float64) # Can't precompute fully
993
994     for j in range(N_users):
995         # Find ratings for user j
996         start_idx = R_train_coo_csc.indptr[j]
997         end_idx = R_train_coo_csc.indptr[j+1]
998         if start_idx == end_idx: # No ratings for this user
999             continue
1000
1001         movie_indices = R_train_coo_csc.indices[start_idx:end_idx]
1002         ratings_centered = R_train_coo_csc.data[start_idx:end_idx]
1003
1004         U_j = U[movie_indices, :] # Movies rated by user j (n_j x R)
1005

```

```

1006     # Adjust ratings by movie bias:  $r_{ij} - \mu - b_i$ 
1007     adjusted_ratings = ratings_centered - movie_bias[movie_indices]
1008
1009     # Calculate  $A = U_j^T U_j + \lambda I$ 
1010     A = U_j.T @ U_j + lambda_sq * np.eye(rank, dtype=np.float64)
1011
1012     # Calculate  $b = U_j^T * \text{adjusted\_ratings}$ 
1013     b = U_j.T @ adjusted_ratings
1014
1015     try:
1016         W[:, j] = np.linalg.solve(A, b)
1017     except np.linalg.LinAlgError:
1018         logger.warning(f"ALS: Solve failed for user {j}, using pseudo-inverse.")
1019         try:
1020             W[:, j] = np.linalg.pinv(A) @ b
1021         except Exception as e_pinv:
1022             logger.error(f"ALS: Pseudo-inverse failed for user {j}: {e_pinv}. Setting W_j to zero.")
1023             W[:, j] = 0.0 # Set to zero vector
1024
1025     return W.astype(np.float64)
1026
1027 def update_movie_factors(
1028     R_train_coo_csr: sparse.csr_matrix, # Centered ratings, CSR format
1029     W: np.ndarray,
1030     user_bias: np.ndarray,
1031     movie_bias: np.ndarray,
1032     lambda_sq: float,
1033     rank: int,
1034     M_movies: int
1035 ) -> np.ndarray:
1036     """Solves for U (movie factors) fixing W and biases."""
1037     N = W.shape[1]
1038     U = np.zeros((M_movies, rank), dtype=np.float64)
1039     # Precompute  $W W^T + \lambda I$  (used in the denominator)
1040     # Note: This is used inside the loop per movie based on specific users  $W_i$ 
1041     #  $W_t W = W @ W.T + \lambda I$  # Can't precompute fully
1042
1043     for i in range(M_movies):
1044         # Find ratings for movie i
1045         start_idx = R_train_coo_csr.indptr[i]
1046         end_idx = R_train_coo_csr.indptr[i+1]
1047         if start_idx == end_idx: # No ratings for this movie
1048             continue
1049
1050         user_indices = R_train_coo_csr.indices[start_idx:end_idx]
1051         ratings_centered = R_train_coo_csr.data[start_idx:end_idx]
1052
1053         W_i = W[:, user_indices] # Users who rated movie i ( $R \times n_i$ )
1054
1055         # Adjust ratings by user bias:  $r_{ij} - \mu - b_u$ 
1056         adjusted_ratings = ratings_centered - user_bias[user_indices]
1057
1058         # Calculate  $A = W_i W_i^T + \lambda I$ 
1059         A = W_i @ W_i.T + lambda_sq * np.eye(rank, dtype=np.float64)
1060
1061         # Calculate  $b = W_i * \text{adjusted\_ratings}$ 
1062         b = W_i @ adjusted_ratings
1063
1064         try:
1065             U[i, :] = np.linalg.solve(A, b)
1066         except np.linalg.LinAlgError:
1067             logger.warning(f"ALS: Solve failed for movie {i}, using pseudo-inverse.")
1068             try:
1069                 U[i, :] = np.linalg.pinv(A) @ b
1070             except Exception as e_pinv:
1071                 logger.error(f"ALS: Pseudo-inverse failed for movie {i}: {e_pinv}. Setting U_i to zero.")
1072                 U[i, :] = 0.0 # Set to zero vector
1073
1074     return U.astype(np.float64)
1075
1076
1077 def update_biases(
1078     R_train_coo: sparse.coo_matrix, # Centered ratings
1079     U: np.ndarray,
1080     W: np.ndarray,
1081     user_bias: np.ndarray,
1082     movie_bias: np.ndarray,
1083     lambda_sq: float
1084 ) -> np.ndarray:
1085     """Solves for user and movie biases"""
1086     # User biases
1087     N = W.shape[1]
1088     user_indices = np.zeros(N, dtype=bool)
1089     for i in range(M_movies):
1090         start_idx = R_train_coo.indptr[i]
1091         end_idx = R_train_coo.indptr[i+1]
1092         if start_idx == end_idx:
1093             continue
1094         user_indices = R_train_coo.indices[start_idx:end_idx]
1095         ratings_centered = R_train_coo.data[start_idx:end_idx]
1096         W_i = W[:, user_indices]
1097         adjusted_ratings = ratings_centered - user_bias[user_indices]
1098         # Calculate  $A = W_i W_i^T + \lambda I$ 
1099         A = W_i @ W_i.T + lambda_sq * np.eye(rank, dtype=np.float64)
1100         # Calculate  $b = W_i * \text{adjusted\_ratings}$ 
1101         b = W_i @ adjusted_ratings
1102         try:
1103             U_i = np.linalg.solve(A, b)
1104         except np.linalg.LinAlgError:
1105             logger.warning(f"ALS: Solve failed for movie {i}, using pseudo-inverse.")
1106             try:
1107                 U_i = np.linalg.pinv(A) @ b
1108             except Exception as e_pinv:
1109                 logger.error(f"ALS: Pseudo-inverse failed for movie {i}: {e_pinv}. Setting U_i to zero.")
1110                 U_i = 0.0 # Set to zero vector
1111         user_indices[user_indices] = U_i
1112
1113     # Movie biases
1114     M_movies = U.shape[0]
1115     movie_indices = np.zeros(M_movies, dtype=bool)
1116     for i in range(M_movies):
1117         start_idx = R_train_coo.indptr[i]
1118         end_idx = R_train_coo.indptr[i+1]
1119         if start_idx == end_idx:
1120             continue
1121         movie_indices = R_train_coo.indices[start_idx:end_idx]
1122         ratings_centered = R_train_coo.data[start_idx:end_idx]
1123         U_i = U[i, :]
1124         adjusted_ratings = ratings_centered - U_i
1125         # Calculate  $A = U_i U_i^T + \lambda I$ 
1126         A = U_i @ U_i.T + lambda_sq * np.eye(rank, dtype=np.float64)
1127         # Calculate  $b = U_i * \text{adjusted\_ratings}$ 
1128         b = U_i @ adjusted_ratings
1129         try:
1130             W_i = np.linalg.solve(A, b)
1131         except np.linalg.LinAlgError:
1132             logger.warning(f"ALS: Solve failed for movie {i}, using pseudo-inverse.")
1133             try:
1134                 W_i = np.linalg.pinv(A) @ b
1135             except Exception as e_pinv:
1136                 logger.error(f"ALS: Pseudo-inverse failed for movie {i}: {e_pinv}. Setting W_i to zero.")
1137                 W_i = 0.0 # Set to zero vector
1138         movie_indices[movie_indices] = W_i
1139
1140     return user_indices, movie_indices

```

```

1083     global_mean: float,
1084     lambda_bias: float,
1085     N_users: int,
1086     M_movies: int
1087 ) -> Tuple[np.ndarray, np.ndarray]:
1088     """Updates user and movie biases based on current residuals."""
1089     new_user_bias = np.zeros_like(user_bias)
1090     new_movie_bias = np.zeros_like(movie_bias)
1091     user_counts = np.zeros_like(user_bias)
1092     movie_counts = np.zeros_like(movie_bias)
1093
1094     # Calculate residuals:  $r_{ij} = \mu - U_i^T W_j$ 
1095     rows, cols, vals_centered = R_train_coo.row, R_train_coo.col, R_train_coo.data
1096     dot_prods = np.array([np.dot(U[r, :], W[:, c]) for r, c in zip(rows, cols)], dtype=np.float64)
1097     residuals = vals_centered - dot_prods # Residual =  $(r_{ij} - \mu) - U_i^T W_j$ 
1098
1099     # Update user biases:  $b_u = \text{sum}(\text{residual} - b_i) / (\text{count} + \text{lambda\_bias})$ 
1100     np.add.at(new_user_bias, cols, residuals - movie_bias[rows])
1101     np.add.at(user_counts, cols, 1)
1102     new_user_bias = new_user_bias / (user_counts + lambda_bias + 1e-9) # Add epsilon for stability
1103
1104     # Update movie biases:  $b_i = \text{sum}(\text{residual} - b_u) / (\text{count} + \text{lambda\_bias})$ 
1105     np.add.at(new_movie_bias, rows, residuals - new_user_bias[cols]) # Use updated user bias
1106     np.add.at(movie_counts, rows, 1)
1107     new_movie_bias = new_movie_bias / (movie_counts + lambda_bias + 1e-9) # Add epsilon for stability
1108
1109     return new_user_bias.astype(np.float64), new_movie_bias.astype(np.float64)
1110
1111 def run_als_with_biases(
1112     R_train_coo: sparse.coo_matrix, # Centered ratings
1113     global_mean: float,
1114     probe_users_mapped: np.ndarray,
1115     probe_movies_mapped: np.ndarray,
1116     probe_ratings_true: np.ndarray,
1117     N_users_active: int,
1118     M_movies_active: int,
1119     rank_local: int,
1120     n_iters: int, # Max iterations
1121     lam_sq: float,
1122     lam_bias: float,
1123     rng: Generator,
1124     init_scale: float = INIT_SCALE_NON_CONVEX,
1125     tol: float = ALS_TOL
1126 ) -> Dict[str, List]:
1127     """Runs Alternating Least Squares with biases."""
1128     logger.info("Starting ALS Solver with Biases...")
1129     U, W, user_bias, movie_bias = initialize_factors_and_biases(
1130         M_movies_active, N_users_active, rank_local, rng, init_scale
1131     )
1132
1133     hist_loss = [] # Loss not typically tracked directly in ALS, focus on RMSE
1134     hist_rmse = []
1135     hist_time = []
1136
1137     start_time = time.time()
1138     last_rmse = np.inf
1139
1140     # Precompute sparse matrix formats for efficiency
1141     R_train_csc = R_train_coo.tocsc()
1142     R_train_csr = R_train_coo.tocsr()
1143
1144     for k_iter in range(1, n_iters + 1):
1145         iter_start_time = time.time()
1146         logger.info(f"--- Starting ALS Iteration {k_iter:02d} ---")
1147
1148         # Update user factors (W)
1149         logger.debug(f"Iter {k_iter}: Updating user factors (W)...")
1150         W = update_user_factors(R_train_csc, U, user_bias, movie_bias, lam_sq, rank_local, N_users_active)
1151
1152         # Update movie factors (U)
1153         logger.debug(f"Iter {k_iter}: Updating movie factors (U)...")
1154         U = update_movie_factors(R_train_csr, W, user_bias, movie_bias, lam_sq, rank_local, M_movies_active)
1155
1156         # Update biases
1157         logger.debug(f"Iter {k_iter}: Updating biases...")
1158         user_bias, movie_bias = update_biases(R_train_coo, U, W, user_bias, movie_bias, global_mean, lam_bias, N_users_active, M_movies)
1159
1160         # Evaluate RMSE

```

```

1161     logger.debug(f"Iter {k_iter}: Evaluating RMSE...")
1162     current_rmse = evaluate_rmse_with_biases(
1163         U, W, user_bias, movie_bias, global_mean,
1164         probe_users_mapped, probe_movies_mapped, probe_ratings_true
1165     )
1166     current_time = time.time() - start_time
1167     hist_rmse.append(current_rmse)
1168     hist_time.append(current_time)
1169
1170     iter_time = time.time() - iter_start_time
1171     logger.info(f"Iter {k_iter:02d}: RMSE = {current_rmse:.6f} (Time: {iter_time:.2f}s)")
1172
1173     # Check convergence
1174     if abs(last_rmse - current_rmse) < tol:
1175         logger.info(f"ALS converged at iteration {k_iter} (RMSE change < {tol})")
1176         break
1177     last_rmse = current_rmse
1178
1179     logger.info("ALS Solver with Biases Finished.")
1180     return {
1181         'loss': [], # ALS doesn't typically track the combined loss easily
1182         'rmse': hist_rmse,
1183         'time': hist_time,
1184         'U': U, 'W': W, 'bu': user_bias, 'bi': movie_bias
1185     }
1186
1187 # --- Stochastic Gradient Single User (NEW - for SARAH/SPIDER) ---
1188 def stochastic_gradient_single_user(U, user_idx, N_users, N_movies, loss_args):
1189     """ Computes the UNSCALED gradient component d L_user_idx / dU for a single user. """
1190     # Unpack loss_args (assumes structure matches loss_and_grad_serial_with_biases)
1191     global_mean, rows_idx, cols_idx, vals_true_centered, _, _, rank_func, lambda_sq_func, lambda_bias_func = loss_args
1192     M, R = U.shape
1193     G_user = np.zeros_like(U, dtype=np.float32)
1194     if user_idx not in user_data_arrays: return G_user # Use precomputed user_data_arrays
1195
1196     W_user_dict = W_closed_efficient(U, N_users, N_movies, user_indices=[user_idx]) # Recompute W for this user
1197     if user_idx not in W_user_dict: return G_user
1198
1199     w_u = W_user_dict[user_idx]
1200     user_data = user_data_arrays[user_idx]
1201     movie_indices = user_data['movies']; rs_t = user_data['rs'] # rs_t are original ratings here
1202     if movie_indices.size == 0: return G_user
1203     if movie_indices.max() >= M or movie_indices.min() < 0: return G_user # Return zero grad if invalid index
1204
1205     # Need centered ratings and biases for gradient calculation
1206     # Recompute biases? Or assume they are passed implicitly? Assume passed via loss_args implicitly (not ideal)
1207     # This function signature needs alignment with how biases are handled if used by SARAH/SPIDER
1208     # For now, approximate using centered ratings and current factors
1209     # This needs refinement if SARAH/SPIDER are primary focus
1210     ratings_centered_user = rs_t - global_mean # Approximate centering
1211
1212     U_k = U[movie_indices, :]
1213     # Need bias terms here for correct error calculation
1214     # Placeholder: Calculate error without biases for now
1215     preds_k_dot = U_k @ w_u
1216     err_k = preds_k_dot - ratings_centered_user # Error against centered rating
1217
1218     grad_vals_k = err_k # Simplified grad without prox term from loss_and_grad
1219     term_k = grad_vals_k.reshape(-1, 1) * w_u.reshape(1, -1)
1220     np.add.at(G_user, movie_indices, term_k.astype(np.float32))
1221     # Add regularization gradient for U rows involved
1222     G_user[movie_indices, :] += lambda_sq_func * U_k
1223
1224     if not np.isfinite(G_user).all():
1225         G_user = np.nan_to_num(G_user, nan=0.0, posinf=0.0, neginf=0.0)
1226     assert G_user.shape == U.shape
1227     return G_user
1228 # --- Euclidean GD Solver (NEW from long.txt, adapted for biases) ---
1229 def run_euclidean_gd(
1230     R_train_coo, global_mean, probe_users_mapped, probe_movies_mapped, probe_ratings_true,
1231     N_users_active, M_movies_active, rank_local, n_iters,
1232     lam_sq, lam_bias, rng, init_scale=INIT_SCALE_NON_CONVEX, lr=1e-7 # Use specific LR
1233 ) -> Dict[str, List]:
1234     """Runs Vanilla Euclidean GD with biases."""
1235     if RANK_MPI == 0: logger.info(f"\n+++ Running Vanilla Euclidean GD (LR={lr:.1e}) +++")
1236     U_euc, W_euc, user_bias, movie_bias = initialize_factors_and_biases(M_movies_active, N_users_active, rank_local, rng, init_scale)
1237     # Note: Euclidean GD doesn't require U to be orthonormal, so we use the direct output
1238     ----

```



```

1238
1239 hist_loss, hist_grad, hist_rmse, hist_time = [], [], [], []; t_start = time.time();
1240 loss_args_biased = (global_mean, R_train_coo.row, R_train_coo.col, R_train_coo.data, M_movies_active, N_users_active, rank_local, l
1241 eval_args_biased = (global_mean, probe_users_mapped, probe_movies_mapped, probe_ratings_true)
1242
1243 try:
1244     current_loss, current_rmse, gU_k, gW_k, gBu_k, gBi_k = record_initial_state_biased(U_euc, W_euc, user_bias, movie_bias, loss_ar
1245     grad_norm_k = np.linalg.norm(gU_k) # Use Euclidean norm for U gradient
1246 except Exception as e:
1247     if RANK_MPI == 0: print(f" ERROR during initial state recording for Euclidean GD: {e}")
1248     return {'loss': [], 'grad_norm': [], 'rmse': [], 'time': []}
1249
1250 if RANK_MPI == 0: hist_loss.append(current_loss); hist_grad.append(grad_norm_k); hist_rmse.append(current_rmse); hist_time.append(t
1251
1252 if RANK_MPI == 0: logger.info("\n Starting Euclidean GD iterations...")
1253 for k in range(n_iters):
1254     iter_t0 = time.time();
1255     # --- inside your Euclidean-GD loop ---
1256     if grad_norm_k < 1e-6:
1257         if RANK_MPI == 0:
1258             logger.info(f"EucGD converged at iter {k}") # or print(...)
1259         break
1260
1261
1262     # Simple Euclidean gradient step for all variables
1263     U_euc -= lr * gU_k
1264     W_euc -= lr * gW_k
1265     user_bias -= lr * gBu_k
1266     movie_bias -= lr * gBi_k
1267
1268     if not (np.isfinite(U_euc).all() and np.isfinite(W_euc).all()):
1269         if RANK_MPI == 0: print(f"EucGD Warning: Non-finite factors at iter {k+1}"); break
1270
1271     try:
1272         current_loss, gU_k, gW_k, gBu_k, gBi_k = loss_and_grad_serial_with_biases(U_euc, W_euc, user_bias, movie_bias, *loss_args_b
1273         current_rmse = evaluate_rmse_with_biases(U_euc, W_euc, user_bias, movie_bias, *eval_args_biased)
1274         grad_norm_k = np.linalg.norm(gU_k) # Euclidean norm
1275         if not (np.isfinite(current_loss) and np.isfinite(gU_k).all() and (np.isnan(current_rmse) or np.isfinite(current_rmse))):
1276             if RANK_MPI == 0: print(f"EucGD Warning: Non-finite values encountered iter {k+1}.")
1277             break
1278     except Exception as e:
1279         if RANK_MPI == 0: print(f"EucGD Error during iteration {k+1}: {e}")
1280         break
1281
1282     if RANK_MPI == 0:
1283         hist_loss.append(current_loss); hist_grad.append(grad_norm_k); hist_rmse.append(current_rmse); hist_time.append(time.time()
1284         if k % 5 == 0 or k == n_iters - 1: print(f" EucGD Iter {k+1:02d} | Loss: {current_loss:.3e} | GradNorm: {grad_norm_k:.3e}
1285
1286 if RANK_MPI == 0: logger.info(f"EucGD finished in {time.time()-t_start:.2f}s");
1287 return {'loss': hist_loss, 'grad_norm': hist_grad, 'rmse': hist_rmse, 'time': hist_time, 'U': U_euc, 'W': W_euc, 'bu': user_bias, '
1288
1289
1290 # ===== #
1291 # CELL 5: Riemannian Solvers (RGD, AGD, Catalyst, DANE) - Renumbered
1292 # ===== #
1293 logger.info("+++ Cell 5: Defining Riemannian Solvers +++")
1294 # --- Stochastic Solvers (SARAH, SPIDER) ---
1295
1296 def run_soft_impute_efficient(
1297     R_train_coo_orig: sparse.coo_matrix, # Original ratings, mapped indices
1298     probe_users_mapped: np.ndarray,
1299     probe_movies_mapped: np.ndarray,
1300     probe_ratings_true: np.ndarray, # Original probe ratings
1301     N_users_active: int,
1302     M_movies_active: int,
1303     n_iters: int,
1304     lambda_reg: float,
1305     k_rank: int, # Initial rank guess / cap for SVD
1306     tol: float,
1307     rng: Generator
1308 ) -> Dict[str, List]:
1309     """ Solves convex problem using efficient Soft-Impute with LinearOperator SVD. """
1310     logger.info("Starting Efficient Convex Soft-Impute Solver (CPU)...")
1311     use_gpu = False # Force CPU as LinearOperator uses SciPy
1312
1313     # Prepare necessary sparse formats of original ratings
1314     R_orig_csr = R_train_coo_orig.tocsr()
1315     R_orig_csc = R_train_coo_orig.tocsc()

```

```

1316 # Create Omega mask (1s where ratings exist)
1317 omega_mask_csr = R_orig_csr.copy(); omega_mask_csr.data[:] = 1
1318 omega_mask_csc = omega_mask_csr.tocsc()
1319
1320 # Initialize factors U, S, V
1321 initial_k = max(1, min(k_rank, M_movies_active, N_users_active))
1322 U = rng.standard_normal(size=(M_movies_active, initial_k)).astype(np.float64) * 0.01
1323 S = np.zeros(initial_k, dtype=np.float64) # Start with S=0 -> Xk=0 initially
1324 V = rng.standard_normal(size=(N_users_active, initial_k)).astype(np.float64) * 0.01
1325 if N_users_active >= initial_k: V, _ = np.linalg.qr(V, mode='reduced') # Orthonormalize V initially
1326
1327 U_old, S_old, V_old = U.copy(), S.copy(), V.copy()
1328 hist_loss, hist_rmse, hist_time, hist_rank = [], [], [], []
1329 start_time = time.time()
1330 current_svd_k = initial_k # Rank for svds call
1331
1332 for k_iter in range(1, n_iters + 1):
1333     iter_start_time = time.time()
1334     logger.info(f"--- Starting SoftImpute Iteration {k_iter:02d} ---")
1335
1336     # Define Linear Operator for Z = P_Omega(R_orig) + P_Omega_Complement(USV^T)
1337     Z_op = ImplicitFillOperator(R_orig_csr, R_orig_csc, omega_mask_csr, omega_mask_csc, U, S, V, (M_movies_active, N_users_active))
1338
1339     # Perform SVD using the LinearOperator
1340     logger.debug(f"Iter {k_iter}: Performing SVD with k={current_svd_k}...")
1341     svd_start_time = time.time()
1342     try:
1343         # Ensure k for svds is valid
1344         k_svds = max(1, min(current_svd_k, M_movies_active - 1, N_users_active - 1))
1345         if k_svds <= 0:
1346             logger.warning(f"Iter {k_iter}: Matrix dimensions too small for SVD. Skipping.")
1347             rank_k = 0; S_new = np.array([], dtype=np.float64)
1348             U_new = np.zeros((M_movies_active, 0), dtype=np.float64)
1349             Vt_new = np.zeros((0, N_users_active), dtype=np.float64) # Need Vt shape
1350         else:
1351             # Use scipy's svds which works with LinearOperator
1352             U_new, S_new_raw, Vt_new = svds(Z_op, k=k_svds, which='LM', tol=1e-4, maxiter=100) # Adjust svds tol/maxiter if needed
1353
1354             # svds returns sorted singular values (largest first) - reverse order
1355             S_new_raw = S_new_raw[::-1]
1356             U_new = U_new[:, ::-1]
1357             Vt_new = Vt_new[::-1, :]
1358
1359             S_new = soft_threshold(S_new_raw, lambda_reg) # Threshold
1360             V_new = Vt_new.T # Transpose Vt to get V
1361             rank_k = int(np.sum(S_new > 1e-10))
1362
1363             logger.debug(f"Iter {k_iter}: SVD finished in {time.time() - svd_start_time:.2f}s. Rank after thresholding: {rank_k}")
1364
1365             if rank_k == 0:
1366                 logger.warning(f"Iter {k_iter}: Rank became zero. Resetting.")
1367                 current_svd_k = 1 # Reset k for next SVD
1368                 U = np.zeros((M_movies_active, 1), dtype=np.float64)
1369                 S = np.zeros(1, dtype=np.float64)
1370                 V = np.zeros((N_users_active, 1), dtype=np.float64)
1371             else:
1372                 U = U_new[:, :rank_k].copy()
1373                 S = S_new[:rank_k].copy()
1374                 V = V_new[:, :rank_k].copy()
1375                 current_svd_k = min(rank_k + 5, CONVEX_RANK_K) # Increase k slightly for next iter, capped
1376
1377     except Exception as e:
1378         logger.error(f"SVD failed during SoftImpute iter {k_iter}: {e}", exc_info=True)
1379         break
1380
1381     # Convergence Check
1382     U_diff_norm = np.linalg.norm(U - U_old, 'fro'); S_diff_norm = np.linalg.norm(S - S_old, 'fro'); V_diff_norm = np.linalg.norm(V - V_old, 'fro')
1383     U_norm = max(1.0, np.linalg.norm(U_old, 'fro')); S_norm = max(1.0, np.linalg.norm(S_old, 'fro')); V_norm = max(1.0, np.linalg.norm(V_old, 'fro'))
1384     relative_diff = max(U_diff_norm / U_norm, S_diff_norm / S_norm, V_diff_norm / V_norm) if U_norm > 0 and S_norm > 0 and V_norm > 0 else 0
1385     logger.debug(f"Iter {k_iter}: Max Rel Factor Diff={relative_diff:.4e}, Rank={rank_k}")
1386
1387     # Evaluate Metrics
1388     eval_start_time = time.time()
1389     try:
1390         # Objective: 0.5 * ||P_Omega(X - R_orig)||_F^2 + lambda * ||X||_*
1391         rows, cols = R_train_coo_orig.row, R_train_coo_orig.col
1392         vals_orig = R_train_coo_orig.data

```

```

1393     preds_at_omega_k = np.array([np.dot(U[r, :], S * V[c, :]) for r, c in zip(rows, cols)], dtype=np.float64)
1394     loss_obs_k = 0.5 * np.sum((preds_at_omega_k - vals_orig)**2)
1395     nuclear_norm_k = np.sum(S)
1396     loss_k = loss_obs_k + lambda_reg * nuclear_norm_k
1397
1398     # RMSE: Predict original scale ratings (USV^T) and compare to true validation ratings
1399     dot_prods_probe = np.array([np.dot(U[m, :], S * V[u, :]) for m, u in zip(probe_movies_mapped, probe_users_mapped)], dtype=r
1400     preds_probe_clamped = np.clip(dot_prods_probe, 1.0, 5.0) # Clamp prediction
1401     valid_true_mask_probe = ~np.isnan(ratings_val_true)
1402     if np.any(valid_true_mask_probe):
1403         mse_probe = np.mean((preds_probe_clamped[valid_true_mask_probe] - ratings_val_true[valid_true_mask_probe])**2)
1404         rmse_k = np.sqrt(mse_probe) if mse_probe >= 0 else np.nan
1405     else: rmse_k = np.nan
1406
1407     except Exception as e: logger.error(f"Error during SoftImpute evaluation: {e}"); loss_k, rmse_k, rank_k = np.nan, np.nan, rank_
1408
1409     eval_time = time.time() - eval_start_time
1410     hist_loss.append(loss_k); hist_rmse.append(rmse_k); hist_time.append(time.time() - start_time); hist_rank.append(rank_k)
1411     U_old, S_old, V_old = U.copy(), S.copy(), V.copy() # Update for next convergence check
1412
1413     iter_time = time.time() - iter_start_time
1414     logger.info(f"Iter {k_iter:02d}: Loss={loss_k:.4e}, RMSE={rmse_k:.4f}, Rank={rank_k}, Rel Diff={relative_diff:.4e} (Eval: {eval
1415
1416     if relative_diff < tol: logger.info(f"Soft-Impute converged at iteration {k_iter}"); break
1417
1418     logger.info("Efficient Convex Soft-Impute Solver Finished.")
1419     return {'loss': hist_loss, 'rmse': hist_rmse, 'time': hist_time, 'rank': hist_rank, 'U': U, 'S': S, 'V': V}
1420
1421
1422 # --- Stochastic Solvers (SARAH, SPIDER) ---
1423 class RiemannianSARAH: # Adapted from long.txt
1424     def __init__(self, R, P, g_i, g_batch, batch_size=100, m=1000, eta=1e-3, rng=None):
1425         self.R, self.P, self.g_i, self.g_batch = R, P, g_i, g_batch
1426         self.B, self.m, self.eta = batch_size, m, eta
1427         self.rng = default_rng(rng) if rng is None else rng
1428     def run(self, U0, n_steps, grad_args, active_idx, sampling_prob=None):
1429         if active_idx is None or len(active_idx) == 0: return U0
1430         rng = self.rng; U = U0.copy().astype(np.float32); v = np.zeros_like(U0, dtype=np.float32)
1431         U_prev = U.copy().astype(np.float32); num_active = len(active_idx)
1432         for t in range(n_steps):
1433             if t % self.m == 0:
1434                 current_batch_size = min(self.B, num_active);
1435                 if current_batch_size == 0: continue
1436                 batch_indices = rng.choice(active_idx, size=current_batch_size, p=sampling_prob, replace=True)
1437                 try: v = self.g_batch(U, batch_indices, *grad_args).astype(np.float32)
1438                 except Exception as e: logger.error(f"SARAH refresh grad error: {e}"); v = np.zeros_like(U)
1439                 if not np.isfinite(v).all(): logger.warning(f"SARAH non-finite refresh grad step {t}"); v = np.zeros_like(U)
1440             else:
1441                 if num_active == 0: continue
1442                 i_idx = rng.choice(active_idx, size=1, p=sampling_prob, replace=True)[0]; i = int(i_idx)
1443                 try:
1444                     v_new = self.g_i(U, i, *grad_args).astype(np.float32)
1445                     v_old = self.g_i(U_prev, i, *grad_args).astype(np.float32)
1446                     if np.isfinite(v_new).all() and np.isfinite(v_old).all(): v += v_new - v_old
1447                 except Exception as e: logger.error(f"SARAH single grad error user {i}: {e}")
1448                 G_proj = self.P(U, v); step = (-self.eta * G_proj).astype(np.float32)
1449                 if should_stop_subproblem(G_proj, step): break
1450                 U_prev = U.copy(); U_next = self.R(U, step)
1451                 if not np.isfinite(U_next).all(): logger.warning(f"SARAH non-finite U step {t+1}"); U = U_prev; break
1452                 U = U_next
1453         return U
1454 class RiemannianSPIDER: # Adapted from long.txt
1455     def __init__(self, retraction, proj, grad_i, grad_batch, m=100, step=1e-3, rng=None):
1456         self.R = retraction; self.P = proj; self.g_i = grad_i; self.g_batch = grad_batch
1457         self.m = m; self.eta = step
1458         self.rng = default_rng(rng) if rng is None else rng
1459     def run(self, U0, n_steps, grad_args, active_idx, sampling_prob=None):
1460         if active_idx is None or len(active_idx) == 0: return U0
1461         rng = self.rng; U = U0.copy().astype(np.float32); v = np.zeros_like(U0, dtype=np.float32)
1462         U_prev = U0.copy().astype(np.float32); num_active = len(active_idx)
1463         for t in range(n_steps):
1464             if t % self.m == 0:
1465                 current_batch_size = min(self.m, num_active); # Use m as batch size for refresh
1466                 if current_batch_size == 0: continue
1467                 batch_indices = rng.choice(active_idx, size=current_batch_size, p=sampling_prob, replace=True)
1468                 try: v = self.g_batch(U, batch_indices, *grad_args).astype(np.float32)
1469                 except Exception as e: logger.error(f"SPIDER refresh grad error: {e}"); v = np.zeros_like(U)
1470                 if not np.isfinite(v).all(): logger.warning(f"SPIDER non-finite refresh grad step {t+1}"); v = np.zeros_like(U)

```

```

1470         if not np.isfinite(v).all(): logger.warning("SPIDER non-finite grad step {t}"); v = np.zeros_like(v)
1471     else:
1472         if num_active == 0: continue
1473         i_idx = rng.choice(active_idx, size=1, p=sampling_prob, replace=True)[0]; i = int(i_idx)
1474         try:
1475             grad_new = self.g_i(U, i, *grad_args).astype(np.float32)
1476             grad_old = self.g_i(U_prev, i, *grad_args).astype(np.float32)
1477             if np.isfinite(grad_new).all() and np.isfinite(grad_old).all(): v = v + grad_new - grad_old
1478         except Exception as e: logger.error(f"SPIDER single grad error user {i}: {e}")
1479         G_proj = self.P(U, v); step_vec = (-self.eta * G_proj).astype(np.float32)
1480         if should_stop_subproblem(G_proj, step_vec): break
1481         U_prev = U.copy(); U_next = self.R(U, step_vec)
1482         if not np.isfinite(U_next).all(): logger.warning(f"SPIDER non-finite U step {t+1}"); U = U_prev; break
1483         U = U_next
1484     return U
1485 # --- RGD Solver ---
1486
1487 def run_rgd_with_biases(
1488     R_train_coo, global_mean, probe_users_mapped, probe_movies_mapped, probe_ratings_true,
1489     N_users_active, M_movies_active, rank_local, n_iters,
1490     lam_sq, lam_bias, rng, init_scale=INIT_SCALE_NON_CONVEX,
1491     lr_init=INIT_LR_RIEMANN, ls_beta=LS_BETA, ls_sigma=LS_SIGMA
1492 ) -> Dict[str, List]:
1493     """Runs Riemannian Gradient Descent with biases."""
1494     logger.info("Starting RGD Solver with Biases...")
1495     U, W, user_bias, movie_bias = initialize_factors_and_biases(M_movies_active, N_users_active, rank_local, rng, init_scale)
1496     hist_loss, hist_rmse, hist_time, hist_grad_norm = [], [], [], []
1497     start_time = time.time(); lr_k = lr_init
1498     loss_args_biased = (global_mean, R_train_coo.row, R_train_coo.col, R_train_coo.data, M_movies_active, N_users_active, rank_local, 1
1499     eval_args_biased = (global_mean, probe_users_mapped, probe_movies_mapped, probe_ratings_true)
1500     try:
1501         loss_k, rmse_k, gU_k, gW_k, gBu_k, gBi_k = record_initial_state_biased(U, W, user_bias, movie_bias, loss_args_biased, eval_args
1502         hist_loss.append(loss_k); hist_rmse.append(rmse_k); hist_time.append(time.time() - start_time)
1503         gU_proj_k = ProjTangent(U, gU_k); hist_grad_norm.append(np.linalg.norm(gU_proj_k))
1504     except Exception as e: logger.error(f"RGD Init Error: {e}"); return {'loss': [], 'rmse': [], 'time': [], 'grad_norm': []}
1505
1506     for k in range(n_iters):
1507         iter_start_time = time.time()
1508         gU_proj_k = ProjTangent(U, gU_k)
1509         grad_norm_k = np.linalg.norm(gU_proj_k)
1510         hist_grad_norm.append(grad_norm_k)
1511
1512         # --- FIX: Check Riemannian Gradient Norm ---
1513         if grad_norm_k < 1e-6: logger.info("RGD Converged (grad norm)"); break
1514         # -----
1515
1516         ls_loss_args = (W, user_bias, movie_bias) + loss_args_biased
1517         lr_step, U_next, loss_next = ArmijoLineSearchRiemannian(U, gU_k, ls_loss_args, loss_k, lr_k, ls_beta, ls_sigma)
1518         if lr_step == 0.0: logger.warning("RGD Line search failed."); break
1519
1520         lr_fixed_other = 1e-4
1521         W -= lr_fixed_other * gW_k; user_bias -= lr_fixed_other * gBu_k; movie_bias -= lr_fixed_other * gBi_k
1522         U = U_next; loss_k = loss_next
1523         lr_k = min(lr_step / np.sqrt(ls_beta), lr_init * 2)
1524
1525         _, gU_k, gW_k, gBu_k, gBi_k = loss_and_grad_serial_with_biases(U, W, user_bias, movie_bias, *loss_args_biased)
1526         rmse_k = evaluate_rmse_with_biases(U, W, user_bias, movie_bias, *eval_args_biased)
1527         hist_loss.append(loss_k); hist_rmse.append(rmse_k); hist_time.append(time.time() - start_time)
1528         iter_time = time.time() - iter_start_time
1529         logger.info(f"Iter {k+1:02d}: Loss={loss_k:.4e}, RMSE={rmse_k:.4f}, GradNorm={grad_norm_k:.2e}, LR={lr_step:.2e} (Time: {iter_t
1530
1531     logger.info("RGD Solver Finished.")
1532
1533 # --- RAGD Solver ---
1534
1535 #
1536 # --- RAGD Solver ---
1537 def run_ragd_with_biases(
1538     R_train_coo, global_mean, probe_users_mapped, probe_movies_mapped, probe_ratings_true,
1539     N_users_active, M_movies_active, rank_local, n_iters,
1540     lam_sq, lam_bias, rng, init_scale=INIT_SCALE_NON_CONVEX,
1541     lr_init=INIT_LR_RIEMANN, ls_beta=LS_BETA, ls_sigma=LS_SIGMA,
1542     gamma=RAGD_GAMMA, mu=RAGD_MU, beta_ragd=RAGD_BETA
1543 ) -> Dict[str, List]:
1544     """Runs Riemannian Accelerated Gradient Descent with biases."""
1545     logger.info("Starting RAGD Solver with Biases...")
1546     U_k, W_k, user_bias_k, movie_bias_k = initialize_factors_and_biases(M_movies_active, N_users_active, rank_local, rng, init_scale)
1547     nu_k = U_k.copy() # Momentum state

```

```

1548 gamma_k = gamma
1549 min_lambda_k = lr_init
1550
1551 hist_loss, hist_rmse, hist_time, hist_grad_norm = [], [], [], []
1552 start_time = time.time()
1553
1554 loss_args_biased = (global_mean, R_train_coo.row, R_train_coo.col, R_train_coo.data, M_movies_active, N_users_active, rank_local, l
1555 eval_args_biased = (global_mean, probe_users_mapped, probe_movies_mapped, probe_ratings_true)
1556
1557 try:
1558     loss_k, rmse_k, gU_k, gW_k, gBu_k, gBi_k = record_initial_state_biased(U_k, W_k, user_bias_k, movie_bias_k, loss_args_biased, e
1559     hist_loss.append(loss_k); hist_rmse.append(rmse_k); hist_time.append(time.time() - start_time)
1560     gU_proj_k = ProjTangent(U_k, gU_k); hist_grad_norm.append(np.linalg.norm(gU_proj_k))
1561 except Exception as e: logger.error(f"RAGD Init Error: {e}"); return {'loss': [], 'rmse': [], 'time': [], 'grad_norm': []}
1562
1563 def solve_alpha_eqn(current_min_lambda, gamma, mu):
1564     a = 1.0; b = current_min_lambda * (gamma - mu); c = -current_min_lambda * gamma
1565     delta = b**2 - 4*a*c
1566     if delta < 0: return 0.0
1567     alpha1 = (-b + np.sqrt(delta))/(2*a); alpha2 = (-b - np.sqrt(delta))/(2*a)
1568     if 0 < alpha1 < 1: return alpha1
1569     if 0 < alpha2 < 1: return alpha2
1570     return 0.0
1571
1572 for k in range(n_iters):
1573     iter_start_time = time.time()
1574     logger.info(f"--- Starting RAGD Iteration {k+1:02d} ---")
1575
1576     alpha = solve_alpha_eqn(min_lambda_k, gamma_k, mu)
1577     if alpha == 0.0: alpha = 1e-6 # Avoid division by zero / stagnation
1578     gamma_bar = (1 - alpha) * gamma_k + alpha * mu
1579     if gamma_bar == 0.0: gamma_bar = 1e-6
1580
1581     # Extrapolation step for y_t (only on U)
1582     logmap_nu_theta = LogMapApprox(U_k, nu_k)
1583     y_t = OrthRetraction(U_k, (alpha * gamma_k / gamma_bar) * logmap_nu_theta)
1584
1585     # Gradient at y_t (need W and biases at y_t? Assume they stay at k for simplicity)
1586     loss_yt, gU_yt, gW_yt, gBu_yt, gBi_yt = loss_and_grad_serial_with_biases(
1587         y_t, W_k, user_bias_k, movie_bias_k, *loss_args_biased
1588     )
1589
1590     # Line search from y_t to find theta_{k+1} (U_{k+1})
1591     ls_loss_args = (W_k, user_bias_k, movie_bias_k) + loss_args_biased
1592     lr_step, U_kp1, loss_kp1 = ArmijoLineSearchRiemannian(
1593         y_t, gU_yt, ls_loss_args, loss_yt, min_lambda_k, ls_beta, ls_sigma
1594     )
1595
1596     if lr_step == 0.0: logger.warning("RAGD Line search failed."); break
1597     min_lambda_k = lr_step # Update min LR found
1598
1599     # Update nu (momentum state)
1600     logmap_nu_yt = LogMapApprox(y_t, nu_k)
1601     grad_proj_yt = ProjTangent(y_t, gU_yt)
1602     nu_update_vec = ((1 - alpha) * gamma_k / gamma_bar) * logmap_nu_yt - (alpha / gamma_bar) * grad_proj_yt
1603     nu_kp1 = OrthRetraction(y_t, nu_update_vec)
1604
1605     # Update W and biases (simple gradient step with decayed LR for stability)
1606     lr_fixed_other = 1e-4 * (0.9**k) # Use a small decaying LR
1607     W_kp1 = W_k - lr_fixed_other * gW_k
1608     user_bias_kp1 = user_bias_k - lr_fixed_other * gBu_k
1609     movie_bias_kp1 = movie_bias_k - lr_fixed_other * gBi_k
1610
1611     # Update state
1612     U_k, W_k, user_bias_k, movie_bias_k = U_kp1, W_kp1, user_bias_kp1, movie_bias_kp1
1613     nu_k = nu_kp1
1614     gamma_k = gamma_bar / (1 + beta_ragd) # Update gamma
1615     loss_k = loss_kp1
1616
1617     # Evaluate and record
1618     rmse_k = evaluate_rmse_with_biases(U_k, W_k, user_bias_k, movie_bias_k, *eval_args_biased)
1619     # Recompute gradient at the final point U_k for norm calculation
1620     _, gU_k_final, gW_k, gBu_k, gBi_k = loss_and_grad_serial_with_biases(U_k, W_k, user_bias_k, movie_bias_k, *loss_args_biased)
1621     gU_proj_k = ProjTangent(U_k, gU_k_final)
1622     grad_norm_k = np.linalg.norm(gU_proj_k)
1623
1624     hist_loss.append(loss_k); hist_rmse.append(rmse_k); hist_time.append(time.time() - start_time)
1625     hist_grad_norm.append(grad_norm_k)

```

```

1625 hist_grad_norm.append(grad_norm_k)
1626
1627 iter_time = time.time() - iter_start_time
1628 logger.info(f"Iter {k+1:02d}: Loss={loss_k:.4e}, RMSE={rmse_k:.4f}, GradNorm={grad_norm_k:.2e}, LR={lr_step:.2e} (Time: {iter_t
1629
1630 if grad_norm_k < 1e-6: logger.info("RAGD Converged (grad norm)"); break
1631
1632 logger.info("RAGD Solver Finished.")
1633 return {'loss': hist_loss, 'rmse': hist_rmse, 'time': hist_time, 'grad_norm': hist_grad_norm, 'U': U_k, 'W': W_k, 'bu': user_bias_k
1634
1635 # --- Catalyst Solver ---
1636 # --- Catalyst Solver (Modified for Stochastic Inner Solvers) ---
1637 def run_catalyst_stochastic( # Renamed from run_catalyst_phi2_with_biases
1638     R_train_coo, global_mean, probe_users_mapped, probe_movies_mapped, probe_ratings_true,
1639     N_users_active, M_movies_active, rank_local, n_iters,
1640     lam_sq, lam_bias, rng, init_scale=INIT_SCALE_NON_CONVEX,
1641     lr_init=INIT_LR_RIEMANN, ls_beta=LS_BETA, ls_sigma=LS_SIGMA,
1642     kappa_0=KAPPA_0, kappa_cvx=KAPPA_CVX, inner_T_epochs=CATALYST_INNER_T_EPOCHS,
1643     inner_S_epochs_base=CATALYST_INNER_S_EPOCHS_BASE,
1644     max_kappa_doublings=MAX_KAPPA_DOUBLINGS,
1645     inner_solver_type=INNER_SOLVER, # NEW: Specify inner solver
1646     inner_solver_lr = RSVRG_LR, # NEW: LR for stochastic inner solver
1647     inner_solver_bs = RSVRG_BATCH_SIZE # NEW: Batch size for stochastic inner solver
1648 ) -> Dict[str, List]:
1649     """Runs Catalyst-Phi2 using a specified stochastic Riemannian solver."""
1650     solver_name = inner_solver_type.upper()
1651     logger.info(f"Starting Catalyst-Phi2 + {solver_name} Solver with Biases...")
1652     theta_k, W_k, user_bias_k, movie_bias_k = initialize_factors_and_biases(M_movies_active, N_users_active, rank_local, rng, init_scal
1653     theta_km1 = theta_k.copy(); tilde_theta_km1 = theta_k.copy()
1654     alpha_k = 1.0; kappa_k = kappa_0
1655     hist_loss, hist_rmse, hist_time, hist_grad_norm = [], [], [], []
1656     phil_grad_hist, phil_dist_hist = [], [] # Rank 0 diagnostics
1657     start_time = time.time()
1658     loss_args_biased = (global_mean, R_train_coo.row, R_train_coo.col, R_train_coo.data, M_movies_active, N_users_active, rank_local, l
1659     eval_args_biased = (global_mean, probe_users_mapped, probe_movies_mapped, probe_ratings_true)
1660     grad_args_stoch = (N_users_active, M_movies_active, loss_args_biased) # Args for stochastic grad funcs
1661     n_data = R_train_coo.nnz # Use number of ratings for epoch length calculation? Or users? Use users.
1662     n_active_users = N_users_active
1663     epoch_len_batches = max(1, n_active_users // inner_solver_bs) if n_active_users > 0 else 1
1664
1665     try:
1666         loss_k, rmse_k, gU_k, gW_k, gBu_k, gBi_k = record_initial_state_biased(theta_k, W_k, user_bias_k, movie_bias_k, loss_args_biase
1667         hist_loss.append(loss_k); hist_rmse.append(rmse_k); hist_time.append(time.time() - start_time)
1668         gU_proj_k = ProjTangent(theta_k, gU_k); hist_grad_norm.append(np.linalg.norm(gU_proj_k))
1669     except Exception as e: logger.error(f"Catalyst-{solver_name} Init Error: {e}"); return {'loss': [], 'rmse': [], 'time': [], 'grad_r
1670
1671 # Instantiate selected inner solver (consistent across ranks)
1672 inner_solver_instance = None
1673 refresh_period_m = max(1, epoch_len_batches // 2) # Example refresh period
1674 solver_args_inner = {
1675     'R': R_fn, 'P': ProjTangent, 'eta': inner_solver_lr,
1676     'g_i': stochastic_gradient_single_user, 'g_batch': stochastic_gradient_batch,
1677     'g_batch': stochastic_gradient_batch, # now resolved 5/4/2025
1678     'rng': default_rng(SEED + 1 + RANK_MPI) # Ensure different RNG streams per rank
1679 }
1680 if inner_solver_type == "sarah": InnerSolverClass = RiemannianSARAH; solver_args_inner.update({'batch_size': inner_solver_bs, 'm':
1681 elif inner_solver_type == "spider": InnerSolverClass = RiemannianSPIDER; solver_args_inner.update({'m': refresh_period_m})
1682 elif inner_solver_type == "svrg": InnerSolverClass = None # SVRG logic remains embedded
1683 else: raise ValueError(f"Unknown INNER_SOLVER: {inner_solver_type}")
1684 if InnerSolverClass: inner_solver_instance = InnerSolverClass(**solver_args_inner)
1685
1686 for k in range(1, n_iters + 1):
1687     iter_start_time = time.time()
1688     logger.info(f"--- Starting Catalyst-{solver_name} Iteration {k:02d} ---")
1689     kappa_step1 = kappa_k; doubling_count = 0
1690     inner_T_steps_budget = epoch_len_batches * inner_T_epochs # Steps budget
1691
1692     logger.debug(f"Iter {k}: Running Phi1 (kappa adaptation)...")
1693     while True:
1694         prox_center = theta_km1.copy()
1695         # --- Run Inner Solver for Step 1 ---
1696         U_inner1 = None
1697         if InnerSolverClass:
1698             try:
1699                 logger.warning(f"Running inner {solver_name} on f, not h_kappa in Phi1.")
1700                 solver_args_run = (grad_args_stoch, unique_users_train, sampling_prob) # Pass active user IDs
1701                 U_inner1 = inner_solver_instance.run(prox_center, inner_T_steps_budget, *solver_args_run)
1702             except Exception as e inner: logger.error(f"Inner {solver name} (Step 1) failed: {e inner}"); U_inner1 = prox_center

```

```

1703 else: # Embedded SVRG for Step 1 subproblem
1704     U_snapshot = prox_center.copy()
1705     G_full_snapshot = np.zeros_like(U_snapshot) # Calculate full gradient estimate
1706     if n_active_users > 0:
1707         num_batches_for_full_grad = max(1, math.ceil(n_active_users / inner_solver_bs / 5))
1708         count_full = 0
1709         for _ in range(num_batches_for_full_grad):
1710             current_batch_size = min(inner_solver_bs, n_active_users)
1711             if current_batch_size == 0: continue
1712             batch_ids_full = GLOBAL_RNG.choice(unique_users_train, size=current_batch_size, p=sampling_prob, replace=True)
1713             try: G_batch = stochastic_gradient_batch(U_snapshot, batch_ids_full, *grad_args_stoch);
1714             except Exception: continue
1715             if np.isfinite(G_batch).all(): G_full_snapshot += G_batch; count_full += 1
1716         if count_full > 0: G_full_snapshot /= count_full
1717     U_inner1_svrg = U_snapshot.copy();
1718     for i_t in range(inner_T_steps_budget):
1719         current_batch_size = min(inner_solver_bs, n_active_users)
1720         if current_batch_size == 0: break
1721         batch_ids = GLOBAL_RNG.choice(unique_users_train, size=current_batch_size, p=sampling_prob, replace=True)
1722         try: g_curr = stochastic_gradient_batch(U_inner1_svrg, batch_ids, *grad_args_stoch); g_ref = stochastic_gradient_b
1723         except Exception: g_curr = np.zeros_like(U_inner1_svrg); g_ref = np.zeros_like(U_inner1_svrg)
1724         if not (np.isfinite(g_curr).all() and np.isfinite(g_ref).all()): continue
1725         G_vr_f = g_curr - g_ref + G_full_snapshot
1726         if REG_DISTANCE == "euclid": G_prox_term = kappa_step1 * (U_inner1_svrg - prox_center);
1727         else: G_prox_term = - kappa_step1 * LogMapApprox(U_inner1_svrg, prox_center)
1728         subprob_G_vr_euclidean = G_vr_f + G_prox_term
1729         G_proj_vr = ProjTangent(U_inner1_svrg, subprob_G_vr_euclidean)
1730         step_vec = (-inner_solver_lr * G_proj_vr).astype(np.float32)
1731         if should_stop_subproblem(G_proj_vr, step_vec): break
1732         U_next_svrg = R_fn(U_inner1_svrg, step_vec)
1733         if not np.isfinite(U_next_svrg).all(): break
1734         U_inner1_svrg = U_next_svrg
1735     U_inner1 = U_inner1_svrg
1736
1737 # --- Check conditions after inner solve ---
1738 theta_bar_k_T = U_inner1;
1739 try: loss_bar_k_T, G_bar_k_T = loss_and_grad_corrected(theta_bar_k_T, *loss_args_biased)
1740 except Exception as e: logger.error(f"Error evaluating bar_theta: {e}"); loss_bar_k_T = np.inf
1741 if not np.isfinite(loss_bar_k_T): kappa_step1 *= 2; doubling_count += 1; continue
1742 conditions_met = False; phi1_grad_norm = np.nan; d_R_approx = np.nan
1743 if RANK_MPI == 0: # Only rank 0 checks conditions
1744     d_R_approx = np.linalg.norm(LogMapApprox(theta_km1, theta_bar_k_T));
1745     h_k_bar = loss_bar_k_T + 0.5 * kappa_step1 * d_R_approx**2;
1746     loss_km1 = hist_loss[-1] if hist_loss else np.inf
1747     descent_cond_met = (h_k_bar <= loss_km1 + 1e-9 * (1 + abs(loss_km1)))
1748     if REG_DISTANCE == "euclid": subprob_grad_bar_k = G_bar_k_T + kappa_step1 * (theta_bar_k_T - theta_km1);
1749     else: subprob_grad_bar_k = G_bar_k_T - kappa_step1 * LogMapApprox(theta_bar_k_T, theta_km1)
1750     proj_grad_h = ProjTangent(theta_bar_k_T, subprob_grad_bar_k)
1751     phi1_grad_norm = np.linalg.norm(proj_grad_h)
1752     stationarity_rhs = kappa_step1 * d_R_approx
1753     stat_cond_met = phi1_grad_norm <= stationarity_rhs + 1e-9 * (1 + stationarity_rhs)
1754     if descent_cond_met and stat_cond_met:
1755         print(f"      Alg phi_1 Conditions MET kappa={kappa_step1:.1e}")
1756         phi1_grad_hist.append(phi1_grad_norm); phi1_dist_hist.append(d_R_approx)
1757         kappa_k_next = update_kappa_adaptive(kappa_step1, phi1_grad_hist, phi1_dist_hist, theta_bar_k_T)
1758         if abs(kappa_k_next - kappa_step1) > 1e-9: print(f"      Adapting kappa next iter: {kappa_step1:.1e} -> {kappa_k_next}
1759         kappa_k = kappa_k_next
1760         conditions_met = True
1761     else: print(f"      Alg phi_1 Conditions NOT MET (Desc:{descent_cond_met}, Stat:{stat_cond_met}) kappa={kappa_step1:.1e}
1762 if COMM and SIZE_MPI > 1: conditions_met = COMM.bcast(conditions_met, root=0); kappa_k = COMM.bcast(kappa_k, root=0) if cor
1763 if conditions_met: break
1764 else:
1765     kappa_step1 *= 2; doubling_count += 1;
1766     if doubling_count >= MAX_KAPPA_DOUBLINGS: logger.warning("Phi1 max kappa doublings reached."); break
1767 if doubling_count >= MAX_KAPPA_DOUBLINGS: logger.error(f"Catalyst Iter {k}: Phi1 failed. Stopping."); break
1768 bar_theta_k = theta_bar_k_T; loss_bar_k = loss_bar_k_T; G_bar_k = G_bar_k_T; kappa_k = kappa_step1
1769 logger.debug(f"Iter {k}: Phi1 finished. Final kappa={kappa_k:.2e}")
1770
1771 # === Step 2: Extrapolation ===
1772 if k == 1: V_extrap_approx = np.zeros_like(theta_km1)
1773 else: V_extrap_approx = LogMapApprox(theta_km1, tilde_theta_km1)
1774 vartheta_k = R_fn(theta_km1, alpha_k * V_extrap_approx);
1775 if not np.isfinite(vartheta_k).all(): logger.error(f"Step 2 non-finite iter {k}. Stopping."); break
1776
1777 # === Step 3: Accelerated Step (using chosen solver) ===
1778 logger.debug(f"Iter {k}: Running Phi2 (accelerated step)...")
1779 prox_center_S = vartheta_k.copy()

```

```

1780 S_k_epochs = math.ceil(inner_S_epochs_base * math.log(k + 1))
1781 max_inner_iter_2 = S_k_epochs * epoch_len_batches
1782
1783 theta_tilde_k = None
1784 if InnerSolverClass:
1785     try:
1786         logger.warning(f"Running inner {solver_name} on f, not h_kappa_cvx in Phi2.")
1787         solver_args_run_S = (grad_args_stoch, unique_users_train, sampling_prob)
1788         theta_tilde_k = inner_solver_instance.run(prox_center_S, max_inner_iter_2, *solver_args_run_S)
1789     except Exception as e_inner_S: logger.error(f"Inner {solver_name} (Step 3) failed: {e_inner_S}"); theta_tilde_k = prox_cer
1790 else: # Embedded SVRG for Step 3 subproblem
1791     U_snapshot_S = prox_center_S
1792     G_full_snapshot_S = np.zeros_like(U_snapshot_S) # Calculate full gradient estimate
1793     if n_active_users > 0:
1794         num_batches_for_full_grad_S = max(1, math.ceil(n_active_users / inner_solver_bs / 5))
1795         count_S_full = 0
1796         for _ in range(num_batches_for_full_grad_S):
1797             current_batch_size_S = min(inner_solver_bs, n_active_users)
1798             if current_batch_size_S == 0: continue
1799             batch_ids_full_S = GLOBAL_RNG.choice(unique_users_train, size=current_batch_size_S, p=sampling_prob, replace=True)
1800             try: G_batch_S = stochastic_gradient_batch(U_snapshot_S, batch_ids_full_S, *grad_args_stoch);
1801             except Exception: continue
1802             if np.isfinite(G_batch_S).all(): G_full_snapshot_S += G_batch_S; count_S_full += 1
1803         if count_S_full > 0: G_full_snapshot_S /= count_S_full
1804     U_inner2_svrg = U_snapshot_S.copy();
1805     for i_s in range(max_inner_iter_2):
1806         current_batch_size_S = min(inner_solver_bs, n_active_users)
1807         if current_batch_size_S == 0: break
1808         batch_ids_S = GLOBAL_RNG.choice(unique_users_train, size=current_batch_size_S, p=sampling_prob, replace=True)
1809         try: g_curr_S = stochastic_gradient_batch(U_inner2_svrg, batch_ids_S, *grad_args_stoch); g_ref_S = stochastic_gradien
1810         except Exception: g_curr_S = np.zeros_like(U_inner2_svrg); g_ref_S = np.zeros_like(U_inner2_svrg)
1811         if not (np.isfinite(g_curr_S).all() and np.isfinite(g_ref_S).all()): continue
1812         G_vr_f_S = g_curr_S - g_ref_S + G_full_snapshot_S
1813         if REG_DISTANCE == "euclid": G_prox_term_S = KAPPA_CVX * (U_inner2_svrg - prox_center_S);
1814         else: G_prox_term_S = - KAPPA_CVX * LogMapApprox(U_inner2_svrg, prox_center_S)
1815         subprob_G_vr_euclidean_S = G_vr_f_S + G_prox_term_S
1816         G_proj_vr_S = ProjTangent(U_inner2_svrg, subprob_G_vr_euclidean_S)
1817         step_vec_S = (-inner_solver_lr * G_proj_vr_S).astype(np.float32)
1818         if should_stop_subproblem(G_proj_vr_S, step_vec_S): break
1819         U_next_S = R_fn(U_inner2_svrg, step_vec_S)
1820         if not np.isfinite(U_next_S).all(): break
1821         U_inner2_svrg = U_next_S
1822     theta_tilde_k = U_inner2_svrg
1823
1824 try: loss_tilde_k, G_tilde_k = loss_and_grad_corrected(theta_tilde_k, *loss_args);
1825 except Exception as e: logger.error(f"Error evaluating tilde_theta: {e}"); loss_tilde_k = np.inf
1826 if not (np.isfinite(loss_tilde_k) and np.isfinite(G_tilde_k).all()): logger.error(f"Step 3 ({solver_name}) failed iter {k}. Stc
1827
1828 # === Step 4, 5, 6 (Consistent) ===
1829 if loss_bar_k <= loss_tilde_k: theta_kp1, loss_kp1, G_kp1, selected = theta_bar_k, loss_bar_k, G_bar_k, "bar"
1830 else: theta_kp1, loss_kp1, G_kp1, selected = theta_tilde_k, loss_tilde_k, G_tilde_k, "tilde"
1831 V_update_approx = LogMapApprox(theta_km1, theta_tilde_k);
1832 tilde_theta_k_next = R_fn(theta_km1, (1.0 / alpha_k) * V_update_approx);
1833 if not np.isfinite(tilde_theta_k_next).all(): logger.error(f"Step 5 non-finite iter {k}. Stopping."); break
1834 alpha_kp1 = (math.sqrt(alpha_k**4 + 4 * alpha_k**2) - alpha_k**2) / 2.0
1835
1836 # --- Update state for next iteration ---
1837 theta_km1 = theta_kp1.copy(); tilde_theta_km1 = tilde_theta_k_next.copy()
1838 alpha_k = alpha_kp1; loss_k = loss_kp1
1839 lr_fixed_other = 1e-4 * (0.9**k)
1840 _, _, gW_kp1, gBu_kp1, gBi_kp1 = loss_and_grad_serial_with_biases(theta_kp1, W_k, user_bias_k, movie_bias_k, *loss_args_biased)
1841 W_k -= lr_fixed_other * gW_kp1; user_bias_k -= lr_fixed_other * gBu_kp1; movie_bias_k -= lr_fixed_other * gBi_kp1
1842
1843 # --- Record History ---
1844 rmse_k = evaluate_rmse_with_biases(theta_kp1, W_k, user_bias_k, movie_bias_k, *eval_args_biased)
1845 gU_proj_k = ProjTangent(theta_kp1, G_kp1); grad_norm_k = np.linalg.norm(gU_proj_k)
1846 hist_loss.append(loss_k); hist_rmse.append(rmse_k); hist_time.append(time.time() - start_time); hist_grad_norm.append(grad_norm
1847 iter_time = time.time() - iter_start_time
1848 logger.info(f"Iter {k:02d}: Loss={loss_k:.4e}, RMSE={rmse_k:.4f}, GradNorm={grad_norm_k:.2e}, Kappa={kappa_k:.2e} (Time: {iter_
1849 if grad_norm_k < 1e-6: logger.info(f"Catalyst-{solver_name} Converged (grad norm)"); break
1850
1851 logger.info(f"Catalyst-{solver_name} Solver Finished.")
1852 if k == n_iters: # Append final grad norm if loop finished normally
1853     _, gU_k_final, _, _, _ = loss_and_grad_serial_with_biases(theta_k, W_k, user_bias_k, movie_bias_k, *loss_args_biased)
1854     gU_proj_k = ProjTangent(theta_k, gU_k_final); hist_grad_norm.append(np.linalg.norm(gU_proj_k))
1855 return {'loss': hist_loss, 'rmse': hist_rmse, 'time': hist_time, 'grad_norm': hist_grad_norm, 'U': theta_k, 'W': W_k, 'bu': user_bi
1856
1857

```



```

1858 # --- DANE Solver ---
1859
1860 # --- DANE Solver ---
1861 def run_dane_with_biases(
1862     R_train_coo, global_mean, probe_users_mapped, probe_movies_mapped, probe_ratings_true,
1863     N_users_active, M_movies_active, rank_local, n_iters,
1864     lam_sq, lam_bias, rng, init_scale=INIT_SCALE_NON_CONVEX,
1865     lr_init=INIT_LR_RIEMANN, ls_beta=LS_BETA, ls_sigma=LS_SIGMA,
1866     kappa=DANE_KAPPA
1867 ) -> Dict[str, List]:
1868     """Runs DANE adaptation with biases."""
1869     logger.info("Starting DANE Solver with Biases...")
1870     theta_k, W_k, user_bias_k, movie_bias_k = initialize_factors_and_biases(M_movies_active, N_users_active, rank_local, rng, init_scal
1871     theta_km1 = theta_k.copy()
1872
1873     hist_loss, hist_rmse, hist_time, hist_grad_norm = [], [], [], []
1874     start_time = time.time()
1875     lr_k = lr_init
1876
1877     loss_args_biased = (global_mean, R_train_coo.row, R_train_coo.col, R_train_coo.data, M_movies_active, N_users_active, rank_local, l
1878     eval_args_biased = (global_mean, probe_users_mapped, probe_movies_mapped, probe_ratings_true)
1879
1880     try:
1881         loss_k, rmse_k, gU_k, gW_k, gBu_k, gBi_k = record_initial_state_biased(theta_k, W_k, user_bias_k, movie_bias_k, loss_args_biase
1882         hist_loss.append(loss_k); hist_rmse.append(rmse_k); hist_time.append(time.time() - start_time)
1883         gU_proj_k = ProjTangent(theta_k, gU_k); hist_grad_norm.append(np.linalg.norm(gU_proj_k))
1884     except Exception as e: logger.error(f"DANE Init Error: {e}"); return {'loss': [], 'rmse': [], 'time': [], 'grad_norm': []}
1885
1886     for k in range(n_iters):
1887         iter_start_time = time.time()
1888         logger.info(f"--- Starting DANE Iteration {k+1:02d} ---")
1889
1890         if k == 0:
1891             grad_combined = gU_k # Use initial gradient for first step
1892         else:
1893             reg_grad = RegularizeGradChordalApprox(theta_k, theta_km1, kappa)
1894             grad_combined = gU_k + reg_grad # gU_k is from end of previous iteration
1895
1896         gU_proj_k = ProjTangent(theta_k, grad_combined)
1897         grad_norm_k = np.linalg.norm(gU_proj_k)
1898         hist_grad_norm.append(grad_norm_k) # Log norm before step
1899
1900         if grad_norm_k < 1e-6: logger.info("DANE Converged (grad norm)"); break
1901
1902         # Line search on U update using combined gradient
1903         ls_loss_args = (W_k, user_bias_k, movie_bias_k) + loss_args_biased
1904         lr_step, U_kp1, loss_kp1 = ArmijoLineSearchRiemannian(
1905             theta_k, grad_combined, ls_loss_args, loss_k, lr_k, ls_beta, ls_sigma
1906         )
1907
1908         if lr_step == 0.0: logger.warning("DANE Line search failed."); break
1909
1910         # Update W and biases (simple gradient step with decayed LR?)
1911         lr_fixed_other = 1e-4 * (0.9**k)
1912         W_kp1 = W_k - lr_fixed_other * gW_k
1913         user_bias_kp1 = user_bias_k - lr_fixed_other * gBu_k
1914         movie_bias_kp1 = movie_bias_k - lr_fixed_other * gBi_k
1915
1916         # Update state
1917         theta_km1 = theta_k.copy() # Store previous U
1918         theta_k = U_kp1
1919         W_k, user_bias_k, movie_bias_k = W_kp1, user_bias_kp1, movie_bias_kp1
1920         loss_k = loss_kp1
1921         lr_k = min(lr_step / np.sqrt(ls_beta), lr_init * 2) # Update LR for next search
1922
1923         # Recompute gradients at new point for next iteration
1924         _, gU_k, gW_k, gBu_k, gBi_k = loss_and_grad_serial_with_biases(theta_k, W_k, user_bias_k, movie_bias_k, *loss_args_biased)
1925         rmse_k = evaluate_rmse_with_biases(theta_k, W_k, user_bias_k, movie_bias_k, *eval_args_biased)
1926
1927         hist_loss.append(loss_k); hist_rmse.append(rmse_k); hist_time.append(time.time() - start_time)
1928
1929         iter_time = time.time() - iter_start_time
1930         logger.info(f"Iter {k+1:02d}: Loss={loss_k:.4e}, RMSE={rmse_k:.4f}, GradNorm={grad_norm_k:.2e}, LR={lr_step:.2e} (Time: {iter_t
1931
1932     logger.info("DANE Solver Finished.")
1933     return {'loss': hist_loss, 'rmse': hist_rmse, 'time': hist_time, 'grad_norm': hist_grad_norm, 'U': theta_k, 'W': W_k, 'bu': user_bi
1934

```

```

1935
1936
1937 # ===== #
1938 # CELL 6: Convex Model Solver (Efficient Soft-Impute) - Renumbered
1939 # ===== #
1940 """
1941 Soft-Impute implementation (Mazumder et al., 2010)
1942 =====
1943 • Works with **NumPy/SciPy** on CPU and **CuPy** on GPU - the backend is
1944   detected automatically.
1945 • Accepts
1946   - `X_incomplete` as a dense `numpy.ndarray` / `cupy.ndarray` *or*
1947     a sparse `scipy.sparse` / `cupyx.scipy.sparse` matrix whose
1948     *missing* entries are encoded as **NaN**.
1949 • Returns either a fully-filled dense array *or* the `(U,S,V)` factors.
1950
1951 This is intentionally self-contained - you can drop the file into any
1952 project (pure Python, no extra deps beyond SciPy/CuPy).
1953 """
1954
1955 from __future__ import annotations
1956
1957 import math
1958 import warnings
1959 from typing import Optional, Tuple, Union
1960
1961 import numpy as _np
1962 from numpy.random import default_rng
1963
1964 try:
1965     import cupy as _cp
1966     import cupyx.scipy.sparse as _cpx_sparse
1967     _HAS_CUPY = True
1968 except ImportError: # GPU unavailable
1969     _cp = None # type: ignore
1970     _HAS_CUPY = False
1971
1972 import scipy.sparse as _sp
1973 from scipy.sparse.linalg import svds as _svds # CPU truncated SVD
1974
1975 Array = Union[_np.ndarray, "_cp.ndarray"] # forward reference for CuPy
1976 Sparse = Union[_sp.spmatrix, "_cpx_sparse.spmatrix"]
1977
1978
1979 # -----
1980 # helpers
1981 # -----
1982
1983 def _to_backend(x: Array | Sparse, use_gpu: bool):
1984     """Move *dense* or *sparse* array to the requested backend."""
1985     if use_gpu and not _HAS_CUPY:
1986         raise RuntimeError("CuPy requested but not installed.")
1987
1988     if use_gpu:
1989         if _HAS_CUPY and isinstance(x, (_cp.ndarray | _cpx_sparse.spmatrix)):
1990             return x # already on GPU
1991         return _cp.asarray(x) if not _sp.issparse(x) else _cpx_sparse.csr_matrix(x)
1992     # -> CPU
1993     if isinstance(x, (_np.ndarray | _sp.spmatrix)):
1994         return x
1995     return _cp.asnumpy(x) if not _sp.issparse(x) else _sp.csr_matrix(x.get())
1996
1997
1998 def _soft_threshold(s: Array, lam: float):
1999     return _np.maximum(s - lam, 0.0)
2000
2001
2002 # -----
2003 # main class
2004 # -----
2005 class SoftImpute:
2006     """Matrix completion via nuclear-norm minimisation.
2007
2008     Parameters
2009     -----
2010     lam : float
2011         Regularisation (shrinkage) parameter  $\lambda$ .
2012     max_rank : int | None, optional

```

```

2012 max_rank : int | None, optional
2013     Maximum rank of the factorisation. Defaults to `min(m, n)`.
2014 max_iters : int, optional
2015     Maximum number of iterations (default 100).
2016 tol : float, optional
2017     Stop when relative change in Frobenius norm < `tol` (default 1e-4).
2018 init_fill_method : {"zero", "mean"}
2019     How to fill missing values in the first iteration.
2020 use_gpu : bool, optional
2021     *True* - try CuPy; *False* - force CPU; *None* - auto-detect.
2022 random_state : int | None
2023     RNG seed for reproducible power-iteration initialisation.
2024 return_factors : bool, default False
2025     If *True* return `(U, S, V)` instead of the filled matrix.
2026 ""
2027
2028 def __init__(
2029     self,
2030     lam: float = 5.0,
2031     *,
2032     max_rank: Optional[int] = None,
2033     max_iters: int = 100,
2034     tol: float = 1e-4,
2035     init_fill_method: str = "zero",
2036     use_gpu: Optional[bool] = None,
2037     random_state: Optional[int] = None,
2038     return_factors: bool = False,
2039 ) -> None:
2040     self.lam = float(lam)
2041     self.max_rank = max_rank
2042     self.max_iters = int(max_iters)
2043     self.tol = float(tol)
2044     if init_fill_method not in {"zero", "mean"}:
2045         raise ValueError("init_fill_method must be 'zero' or 'mean'")
2046     self.init_fill_method = init_fill_method
2047     self.use_gpu = (_HAS_CUPY if use_gpu is None else bool(use_gpu))
2048     self.rng = default_rng(random_state)
2049     self.return_factors = return_factors
2050
2051     # will be initialised in `fit_transform`
2052     self.U_: Optional[Array] = None
2053     self.S_: Optional[Array] = None
2054     self.V_: Optional[Array] = None
2055
2056 # -----
2057 def fit_transform(self, X: Array | Sparse) -> Array | Tuple[Array, Array, Array]:
2058     """Run Soft-Impute and return the completed matrix or the factors."""
2059
2060     # move data to desired backend
2061     X = _to_backend(X, self.use_gpu)
2062     xp = _cp if (self.use_gpu) else _np
2063     spmod = _cpx_sparse if (self.use_gpu) else _sp
2064
2065     # sparse -> dense with NaNs where missing -----
2066     if spmod.issparse(X):
2067         X = X.tocsr()
2068         m, n = X.shape
2069         dense = xp.full((m, n), xp.nan, dtype=xp.float32)
2070         rows, cols = X.nonzero()
2071         dense[rows, cols] = X.data.astype(xp.float32)
2072         X = dense
2073     else:
2074         X = X.astype(xp.float32)
2075
2076     nan_mask = xp.isnan(X)
2077     m, n = X.shape
2078     max_rank = self.max_rank or min(m, n)
2079
2080     # initial fill -----
2081     X_filled = X.copy()
2082     if self.init_fill_method == "mean":
2083         col_means = xp.nanmean(X, axis=0)
2084         inds = nan_mask
2085         X_filled[inds] = col_means[xp.newaxis, :][inds]
2086     else: # zero
2087         X_filled[nan_mask] = 0.0
2088
2089     # main iteration -----

```

```

2090     prev_norm = xp.linalg.norm(X_filled)
2091     for it in range(1, self.max_iters + 1):
2092         # truncated SVD: cpu → scipy.sparse.linalg.svds; gpu → full svd of cuPy
2093         if self.use_gpu:
2094             U, S, Vt = xp.linalg.svd(X_filled, full_matrices=False)
2095             U, S, Vt = U[:, :max_rank], S[:max_rank], Vt[:max_rank, :]
2096         else:
2097             # work with float64 for SciPy stability
2098             U, S, Vt = _svds(_sp.csr_matrix(X_filled), k=max_rank, which="LM")
2099             # SciPy returns in ascending order
2100             U, S, Vt = U[:, ::-1], S[::-1], Vt[::-1, :]
2101
2102         # soft-threshold singular values -----
2103         S_shrink = _soft_threshold(S, self.lam)
2104         rank_k = int((S_shrink > 0).sum())
2105         if rank_k == 0:
2106             warnings.warn("All singular values shrunk to 0 - returning previous iterate.")
2107             break
2108         U = U[:, :rank_k]
2109         S_shrink = S_shrink[:rank_k]
2110         Vt = Vt[:rank_k, :]
2111
2112         # reconstruct and impute -----
2113         X_hat = (U * S_shrink) @ Vt      # U (m×r) * diag(S) * V^T (r×n)
2114         X_filled[nan_mask] = X_hat[nan_mask]
2115
2116         # convergence check -----
2117         frob_norm = xp.linalg.norm(X_filled)
2118         rel_change = xp.linalg.norm(X_filled - X_hat) / max(1.0, frob_norm)
2119         if rel_change < self.tol:
2120             break
2121         prev_norm = frob_norm
2122
2123         # store factors on CPU for compat -----
2124         self.U_ = _cp.asnumpy(U) if self.use_gpu else U
2125         self.S_ = _cp.asnumpy(S_shrink) if self.use_gpu else S_shrink
2126         self.V_ = _cp.asnumpy(Vt.T) if self.use_gpu else Vt.T
2127
2128         if self.return_factors:
2129             return self.U_, self.S_, self.V_
2130         return _cp.asnumpy(X_filled) if self.use_gpu else X_filled
2131
2132     # -----
2133     def transform(self, X_new: Array | Sparse) -> Array:
2134         """Impute a *new* matrix with the learnt factors (no retraining)."""
2135         if self.U_ is None:
2136             raise RuntimeError("call fit_transform first")
2137         X_new = _to_backend(X_new, self.use_gpu)
2138         xp = _cp if self.use_gpu else _np
2139         dense = X_new.copy()
2140         nan_mask = xp.isnan(dense)
2141         X_hat = (self.U_ * self.S_) @ self.V_.T
2142         dense[nan_mask] = X_hat[nan_mask]
2143         return _cp.asnumpy(dense) if self.use_gpu else dense
2144     # ===== #
2145     # CELL 7: Run Solvers and Compare Results - Renumbered
2146     # ===== #
2147     logger.info("+++ Cell 7: Running Solvers and Comparing Results +++")
2148
2149     all_results = {}
2150     # --- Initialize Trajectory Cache (Rank 0 only) ---
2151     TRAJECTORY_CACHE = [] if RANK_MPI == 0 else None
2152
2153     # --- Update solver_args with new variable names ---
2154     solver_args = {
2155         "R_train_coo": R_train_coo, "global_mean": global_mean_rating,
2156         "probe_users_mapped": user_ids_val_final, "probe_movies_mapped": movie_ids_val_final,
2157         "probe_ratings_true": ratings_val_true, "N_users_active": N_users_active,
2158         "M_movies_active": M_movies_active, "rank_local": RANK, "lam_sq": LAM_SQ,
2159         "lam_bias": LAM_BIAS, "rng": GLOBAL_RNG, "init_scale": INIT_SCALE_NON_CONVEX,
2160     }
2161
2162     # --- Run Non-Convex Solvers ---
2163     if DATA_AVAILABLE and R_train_coo.nnz > 0 and N_users_active > 0 and M_movies_active > 0:
2164         # Euclidean GD (NEW)
2165         if RANK_MPI == 0: logger.info("\n--- Running Non-Convex Solver (Euclidean GD with Biases) ---")
2166         try: all_results['Non-Convex (EucGD+Bias)'] = run_euclidean_gd(**solver_args, n_iters=N_ITERS_ALL, lr=1e-7) # Added call, specify L

```

```

2167 except Exception as e: logger.error(f"SVRG Failed: {e}", exc_info=True); all_results['Non-Convex (SVRG+Bias)'] = {}
2168 # SVRG
2169 if RANK_MPI == 0: logger.info("\n--- Running Non-Convex Solver (SVRG Adaptation with Biases) ---")
2170 try: all_results['Non-Convex (SVRG+Bias)'] = run_non_convex_svrg_with_biases(**solver_args, n_epochs=N_ITERS_ALL, inner_lr=INIT_LR
2171 except Exception as e: logger.error(f"SVRG Failed: {e}", exc_info=True); all_results['Non-Convex (SVRG+Bias)'] = {}
2172 # ALS
2173 if RANK_MPI == 0: logger.info("\n--- Running Non-Convex Solver (ALS with Biases) ---")
2174 try: all_results['Non-Convex (ALS+Bias)'] = run_als_with_biases(**solver_args, n_iters=N_ITERS_ALL, tol=ALS_TOL)
2175 except Exception as e: logger.error(f"ALS Failed: {e}", exc_info=True); all_results['Non-Convex (ALS+Bias)'] = {}
2176 # RGD
2177 if RANK_MPI == 0: logger.info("\n--- Running Non-Convex Solver (RGD with Biases) ---")
2178 try: all_results['Non-Convex (RGD+Bias)'] = run_rgd_with_biases(**solver_args, n_iters=N_ITERS_ALL, lr_init=INIT_LR_RIEMANN, ls_bet
2179 except Exception as e: logger.error(f"RGD Failed: {e}", exc_info=True); all_results['Non-Convex (RGD+Bias)'] = {}
2180 # RAGD
2181 if RANK_MPI == 0: logger.info("\n--- Running Non-Convex Solver (RAGD with Biases) ---")
2182 try: all_results['Non-Convex (RAGD+Bias)'] = run_ragd_with_biases(**solver_args, n_iters=N_ITERS_ALL, lr_init=INIT_LR_RIEMANN, ls_b
2183 except Exception as e: logger.error(f"RAGD Failed: {e}", exc_info=True); all_results['Non-Convex (RAGD+Bias)'] = {}
2184 # Catalyst + Selected Inner Solver
2185 if RANK_MPI == 0: logger.info(f"\n--- Running Non-Convex Solver (Catalyst-{INNER_SOLVER.upper()} with Biases) ---")
2186 try: all_results[f'Non-Convex (Catalyst+{INNER_SOLVER.upper()})'] = run_catalyst_stochastic(**solver_args, n_iters=N_ITERS_ALL, lr_
2187 except Exception as e: logger.error(f"Catalyst-{INNER_SOLVER.upper()} Failed: {e}", exc_info=True); all_results[f'Non-Convex (Catal
2188 # DANE
2189 if RANK_MPI == 0: logger.info("\n--- Running Non-Convex Solver (DANE with Biases) ---")
2190 try: all_results['Non-Convex (DANE+Bias)'] = run_dane_with_biases(**solver_args, n_iters=N_ITERS_ALL, lr_init=INIT_LR_RIEMANN, ls_b
2191 except Exception as e: logger.error(f"DANE Failed: {e}", exc_info=True); all_results['Non-Convex (DANE+Bias)'] = {}
2192 else:
2193     if RANK_MPI == 0: logger.warning("Skipping Non-Convex Solvers due to missing data or zero dimensions.")
2194
2195 # --- Run Convex Solver (Efficient Soft-Impute) ---
2196 if DATA_AVAILABLE and R_train_coo_orig.nnz > 0 and N_users_active > 0 and M_movies_active > 0:
2197     if RANK_MPI == 0: logger.info("\n--- Running Convex Solver (Efficient Soft-Impute) ---")
2198     try:
2199         results_convex = run_soft_impute_efficient(
2200             R_train_coo_orig=R_train_coo_orig, # Use original ratings matrix
2201             probe_users_mapped=user_ids_val_final,
2202             probe_movies_mapped=movie_ids_val_final,
2203             probe_ratings_true=ratings_val_true, # Use validation ratings
2204             N_users_active=N_users_active,
2205             M_movies_active=M_movies_active,
2206             n_iters=N_ITERS_ALL, # Use N_ITERS_ALL for consistency
2207             lambda_reg=LAM, # Use LAM directly
2208             k_rank = CONVEX_RANK_K,
2209             tol=SOFT_IMPUTE_TOL,
2210             rng=GLOBAL_RNG
2211         )
2212         all_results['Convex (SoftImpute Eff.)'] = results_convex
2213     except Exception as e:
2214         logger.error(f"Failed to run Efficient Soft-Impute Solver: {e}", exc_info=True)
2215         all_results['Convex (SoftImpute Eff.)'] = {'loss': [], 'rmse': [], 'time': [], 'rank': []}
2216 else:
2217     if RANK_MPI == 0: logger.warning("Skipping Convex Solver due to missing data or zero dimensions.")
2218     all_results['Convex (SoftImpute Eff.)'] = {'loss': [], 'rmse': [], 'time': [], 'rank': []}
2219
2220
2221 # --- Plotting Comparison ---
2222 if RANK_MPI == 0:
2223     logger.info("\n--- Generating Comparison Plots ---")
2224     plt.style.use('seaborn-v0_8-whitegrid')
2225     fig, axes = plt.subplots(3, 2, figsize=(12, 11), sharex='col')
2226     fig.suptitle(
2227         f'MovieLens 1M ({RATING_LIMIT/1e6 if RATING_LIMIT else "Full"} M ratings subset), '
2228         f'Rank={RANK}, Outer iters={N_ITERS_ALL}',
2229         fontsize=14,
2230     )
2231
2232 # ----- style dictionary (matches earlier section) -----
2233 styles = {
2234     'Non-Convex (SVRG+Bias)': dict(label=r'SVRG+Bias', style='-', 'p'), alpha=.90, color='tab:purple'),
2235     'Non-Convex (ALS+Bias)': dict(label=r'ALS+Bias', style='-', 'v'), alpha=.90, color='tab:brown'),
2236     'Non-Convex (RGD+Bias)': dict(label=r'RGD+Bias', style='--', 'o'), alpha=.80, color='tab:blue'),
2237     'Non-Convex (RAGD+Bias)': dict(label=r'RAGD+Bias', style='-', 'D'), alpha=.80, color='tab:orange'),
2238     f'Non-Convex (Catalyst+{INNER_SOLVER.upper()})': dict(label=f'Catalyst+{INNER_SOLVER.upper()}', style='-', 's'), alpha=.90, c
2239     'Non-Convex (DANE+Bias)': dict(label=r'DANE+Bias', style='-', 'x'), alpha=.80, color='tab:cyan'),
2240     'Non-Convex (EucGD+Bias)': dict(label=r'EucGD+Bias', style=':', '^'), alpha=.70, color='tab:green'),
2241     'Convex (SoftImpute Eff.)': dict(label=r'SoftImpute (Eff)', style='-', '*'), alpha=.90, color='tab:pink'),
2242 }
2243
2244 # ----- helper for plotting one method -----

```

```

2245 def _plot(ax_iter, ax_time, data, meta):
2246     ls, mk = meta['style']
2247     kw = dict(linestyle=ls, marker=mk, markersize=3, alpha=meta['alpha'], color=meta.get('color', None))
2248     n_loss = len(data.get('loss', [])); n_grad = len(data.get('grad_norm', [])); n_rmse = len(data.get('rmse', [])); n_time = len(c
2249     n = min(n_loss if n_loss > 0 else float('inf'), n_grad if n_grad > 0 else float('inf'), n_rmse if n_rmse > 0 else float('inf'),
2250     if n == float('inf') or n < 2: logger.warning(f" • insufficient points for {meta['label']}"); return
2251
2252     it = np.arange(n)
2253     loss_vals = np.array(data.get('loss', [np.nan]*n)[:n]); grad_vals = np.array(data.get('grad_norm', [np.nan]*n)[:n])
2254     rmse_vals = np.array(data.get('rmse', [np.nan]*n)[:n]); time_vals = np.array(data.get('time', [np.nan]*n)[:n])
2255
2256     # Determine primary metric for grad plot (grad_norm, or gU_norm for SVRG)
2257     grad_metric = grad_vals
2258     if not np.any(np.isfinite(grad_metric)) and 'gU_norm' in data:
2259         grad_metric = np.array(data.get('gU_norm', [np.nan]*n)[:n])
2260
2261     loss_ok = np.isfinite(loss_vals); grad_ok = np.isfinite(grad_metric); rmse_ok = np.isfinite(rmse_vals); time_ok = np.isfinite(t
2262
2263     # iteration domain
2264     if np.any(loss_ok): ax_iter[0].semilogy(it[loss_ok], loss_vals[loss_ok], label=meta['label'], **kw)
2265     if np.any(grad_ok): ax_iter[1].semilogy(it[grad_ok], grad_metric[grad_ok], **kw)
2266     if np.any(rmse_ok): ax_iter[2].plot(it[rmse_ok], rmse_vals[rmse_ok], **kw)
2267
2268     # wall-clock domain
2269     if np.any(loss_ok & time_ok): ax_time[0].semilogy(time_vals[loss_ok & time_ok], loss_vals[loss_ok & time_ok], **kw)
2270     if np.any(grad_ok & time_ok): ax_time[1].semilogy(time_vals[grad_ok & time_ok], grad_metric[grad_ok & time_ok], **kw)
2271     if np.any(rmse_ok & time_ok): ax_time[2].plot(time_vals[rmse_ok & time_ok], rmse_vals[rmse_ok & time_ok], **kw)
2272
2273     # ----- draw every available method -----
2274     for m, d in all_results.items():
2275         if m in styles and d: # Check if history dict is not empty
2276             _plot(axes[:, 0], axes[:, 1], d, styles[m])
2277         else:
2278             logger.warning(f" • no style or no results for '{m}', skipped.")
2279
2280     # labels / titles
2281     axes[0,0].set_ylabel('Objective'); axes[0,0].set_title('Loss vs Iterations')
2282     axes[1,0].set_ylabel(r'$\|\nabla\|$'); axes[1,0].set_title('Grad-norm vs Iterations')
2283     axes[2,0].set_ylabel('Validation RMSE'); axes[2,0].set_xlabel('Iteration k'); axes[2,0].set_title('RMSE vs Iterations')
2284     axes[0,1].set_xscale('log'); axes[0,1].set_ylabel('Objective'); axes[0,1].set_title('Loss vs Wall-time')
2285     axes[1,1].set_xscale('log'); axes[1,1].set_ylabel(r'$\|\nabla\|$'); axes[1,1].set_title('Grad-norm vs Wall-time')
2286     axes[2,1].set_xscale('log'); axes[2,1].set_ylabel('Validation RMSE'); axes[2,1].set_xlabel('Seconds'); axes[2,1].set_title('RMSE vs
2287
2288     for ax in axes.flatten():
2289         ax.grid(True, which='both', linestyle=':', linewidth=.5)
2290         handles, labels = ax.get_legend_handles_labels()
2291         if handles: ax.legend() # Only add legend if there are labeled artists
2292
2293     plt.tight_layout(rect=[0, 0.03, 1, 0.95])
2294     plt.show()
2295
2296     # ----- optional PCA trajectory plot -----
2297     if PCA_AVAILABLE and TRAJECTORY_CACHE is not None and len(TRAJECTORY_CACHE) >= 3:
2298         logger.info("\n+++ Generating PCA Trajectory Plot +++")
2299         try:
2300             traj_dim = TRAJECTORY_CACHE[0].size
2301             valid_traj = [t for t in TRAJECTORY_CACHE if isinstance(t, np.ndarray) and t.size == traj_dim]
2302             if len(valid_traj) >= 3:
2303                 pcs = PCA(n_components=2).fit_transform(np.vstack(valid_traj))
2304                 plt.figure(figsize=(4.5,4)); plt.plot(pcs[:,0], pcs[:,1], '-o', markersize=3)
2305                 plt.title('Optimisation Trajectory (PCA)'); plt.xlabel('PC1'); plt.ylabel('PC2')
2306                 plt.tight_layout(); plt.show()
2307             else: logger.warning("Not enough valid trajectory points for PCA plot.")
2308         except Exception as e_pca: logger.error(f"PCA Trajectory plot failed: {e_pca}")
2309
2310
2311 # --- Final Summary Table ---
2312 if RANK_MPI == 0:
2313     logger.info("\n--- Final Comparison Summary ---")
2314     print(f"{'Method':<30} | {'Final RMSE':<15} | {'Final Loss':<15} | {'Final Rank/GradNorm':<18} | {'Time (s)':<15}")
2315     print(f"{'-'*30}-|{'-'*15}-|{'-'*15}-|{'-'*18}-|{'-'*15}")
2316     def get_last_finite(history, key):
2317         if not isinstance(history, dict): return np.nan
2318         data = history.get(key)
2319         if isinstance(data, (list, np.ndarray)) and len(data) > 0:
2320             arr = np.array(data); finite_vals = arr[np.isfinite(arr)]
2321             return finite_vals[-1] if finite_vals.size > 0 else np.nan

```

```

2322     return np.nan
2323 for label, history in all_results.items():
2324     if not history: print(f"{label:<30} | {'FAILED':<15} | {'FAILED':<15} | {'N/A':<18} | {'N/A':<15}"); continue
2325     final_rmse = get_last_finite(history, 'rmse')
2326     final_loss = get_last_finite(history, 'loss')
2327     final_time = get_last_finite(history, 'time')
2328     final_rank = get_last_finite(history, 'rank') if 'rank' in history else RANK
2329     final_grad_norm = get_last_finite(history, 'grad_norm') if 'grad_norm' in history else np.nan
2330     final_gU_norm = get_last_finite(history, 'gU_norm') if 'gU_norm' in history else np.nan
2331     rmse_str = f"{final_rmse:.6f}" if np.isfinite(final_rmse) else 'NaN'
2332     loss_str = f"{final_loss:.6e}" if np.isfinite(final_loss) and 'ALS' not in label and 'SoftImpute' not in label else 'N/A'
2333     rank_or_grad_str = 'N/A'
2334     if 'SoftImpute' in label: rank_or_grad_str = f"Rank={int(final_rank)}" if np.isfinite(final_rank) else 'N/A'
2335     elif 'grad_norm' in history and np.isfinite(final_grad_norm): rank_or_grad_str = f"||G||={final_grad_norm:.2e}"
2336     elif 'gU_norm' in history and np.isfinite(final_gU_norm): rank_or_grad_str = f"||gU||={final_gU_norm:.2e}"
2337     else: rank_or_grad_str = f"Rank={RANK}"
2338     time_str = f"{final_time:.4f}" if np.isfinite(final_time) else 'N/A'
2339     print(f"{label:<30} | {rmse_str:<15} | {loss_str:<15} | {rank_or_grad_str:<18} | {time_str:<15}")
2340 print("\nComparison Complete.")
2341
2342 # --- ADDED Block 6-a: Run OT Demo (Rank 0 only) ---
2343 # --- ADDED Block 6-a: Run OT Demo (Rank 0 only) ---
2344 if RANK_MPI == 0 and OT_AVAILABLE:
2345     logger.info("\n+++ Running OT Barycentre Demo +++")
2346     try:
2347         ot_demo_results = run_barycentre_demo()
2348         # Optionally plot or process ot_demo_results
2349         plt.figure(figsize=(6, 4))
2350         plt.plot(ot_demo_results['grid'], ot_demo_results['sources'], '--', label='Sources')
2351         plt.plot(ot_demo_results['grid'], ot_demo_results['barycenter'], 'r-', label='Barycenter')
2352         plt.title('Wasserstein Barycenter Demo')
2353         plt.legend(); plt.tight_layout(); plt.show()
2354     except Exception as e_ot:
2355         logger.error(f"OT Barycentre Demo failed: {e_ot}")
2356
2357 # === ADDED Block 6: PCA Trajectory Plot (Rank 0 only) ===
2358 if RANK_MPI == 0 and PCA_AVAILABLE and len(TRAJECTORY_CACHE) >= 3:
2359     logger.info("\n+++ Generating PCA Trajectory Plot +++")
2360     try:
2361         # Ensure all trajectories have the same dimension (flattened U)
2362         traj_dim = TRAJECTORY_CACHE[0].size
2363         valid_traj = [t for t in TRAJECTORY_CACHE if t.size == traj_dim]
2364         if len(valid_traj) >= 3:
2365             pcs = PCA(n_components=2).fit_transform(np.vstack(valid_traj))
2366             plt.figure(figsize=(4.5,4)); plt.plot(pcs[:,0], pcs[:,1], '-o', markersize=3)
2367             plt.title('Optimisation Trajectory (PCA)'); plt.xlabel('PC1'); plt.ylabel('PC2')
2368             plt.tight_layout(); plt.show()
2369         else:
2370             logger.warning("Not enough valid trajectory points for PCA plot.")
2371     except Exception as e_pca:
2372         logger.error(f"PCA Trajectory plot failed: {e_pca}")
2373
2374 # === ADDED Block 7: Dump TeX skeleton to Drive (Rank 0 only) ===
2375 if RANK_MPI == 0:
2376     TEX_PATH = Path(DATA_DIR_STR) / "proofs.tex" # Use Path object
2377     if TEX_PATH.parent.is_dir():
2378         logger.info(f"\n+++ Checking/Writing TeX Proof Skeleton to: {TEX_PATH} +++")
2379         if not TEX_PATH.exists():
2380             try:
2381                 with open(TEX_PATH, "w") as f: f.write(r"\"...\"") # TeX content omitted for brevity
2382                 logger.info(f" Wrote TeX scaffold to {TEX_PATH}")
2383                 except IOError as e: logger.error(f" Error writing TeX file: {e}")
2384             else: logger.info(f" TeX scaffold already exists at {TEX_PATH}, not overwritten.")
2385         else: logger.warning(f" Parent directory for TeX file not found: {TEX_PATH.parent}")
2386
2387
2388 # --- Mount Drive (if not already mounted) ---
2389 if RANK_MPI == 0 and not Path("/content/drive").is_mount():
2390     try:
2391         from google.colab import drive
2392         drive.mount("/content/drive", force_remount=False)
2393     except Exception as e:
2394         logger.error(f"Failed to mount Google Drive at the end: {e}")
2395
2396 # ===== #
2397 # CELL 8: Plots & Dashboards (from long.txt) - Renumbered
2398 # ===== #
2399 if RANK_MDT == 0:

```

```

2400 logger.info("\n+++ Cell 8: Plots & Dashboards +++")
2401
2402 # ----- helper ----- #
2403 def _plot_metric(metric_key: str,
2404                 ylabel: str,
2405                 x_key: str = "time",
2406                 title: str | None = None,
2407                 logy: bool = False,
2408                 logx: bool = True,          # Default: log time axis
2409                 figsize=(8, 5)) -> None:
2410     plt.figure(figsize=figsize)
2411     has_data_to_plot = False
2412
2413     # style dictionary -----
2414     styles = {
2415         'Non-Convex (SVRG+Bias)': dict(label='SVRG+Bias', style='-', 'p'), alpha=.90, color='tab:purple'),
2416         'Non-Convex (ALS+Bias)': dict(label='ALS+Bias', style='-', 'v'), alpha=.90, color='tab:brown'),
2417         'Non-Convex (RGD+Bias)': dict(label='RGD+Bias', style='--', 'o'), alpha=.80, color='tab:blue'),
2418         'Non-Convex (RAGD+Bias)': dict(label='RAGD+Bias', style='-.', 'D'), alpha=.80, color='tab:orange'),
2419         f'Non-Convex (Catalyst+{INNER_SOLVER.upper()})':
2420             dict(label=f'Catalyst+{INNER_SOLVER.upper()}',
2421                 style='-', 's'), alpha=.90, color='tab:red'),
2422         'Non-Convex (DANE+Bias)': dict(label='DANE+Bias', style='-', 'x'), alpha=.80, color='tab:cyan'),
2423         'Non-Convex (EucGD+Bias)': dict(label='EucGD+Bias', style=':', '^'), alpha=.70, color='tab:green'),
2424         'Convex (SoftImpute Eff.)': dict(label='SoftImpute (Eff)', style='-', '*'), alpha=.90, color='tab:pink'),
2425     }
2426
2427     # loop over solver results -----
2428     for name, res in all_results.items():
2429         y = res.get(metric_key, [])
2430         x = res.get(x_key, list(range(len(y)))) if x_key else list(range(len(y)))
2431
2432         if len(y) == 0:
2433             continue
2434
2435         x = np.asarray(x, dtype=float)
2436         y = np.asarray(y, dtype=float)
2437         valid = np.isfinite(x) & np.isfinite(y)
2438         x_plot, y_plot = x[valid], y[valid]
2439
2440         if x_plot.size == 0:
2441             logger.warning(f"No finite data to plot for {name} - {metric_key}")
2442             continue
2443
2444         style = styles.get(name, {})
2445         plt.plot(
2446             x_plot, y_plot,
2447             linestyle=style.get('style', ('-', 'o'))[0],
2448             marker=style.get('style', ('-', 'o'))[1],
2449             markersize=3,
2450             alpha=style.get('alpha', 0.8),
2451             color=style.get('color'),
2452             label=style.get('label', name)
2453         )
2454         has_data_to_plot = True
2455
2456     # axes / formatting -----
2457     plt.xlabel("wall-clock (s)" if x_key == "time" else "iteration")
2458     plt.ylabel(ylabel)
2459     if logx:
2460         plt.xscale("log")
2461     if logy:
2462         plt.yscale("log")
2463     plt.title(title or f"{ylabel} vs {'time' if x_key == 'time' else 'iteration'}")
2464     if has_data_to_plot:
2465         plt.legend()
2466     plt.grid(alpha=.3, which='both', linestyle=':')
2467     plt.tight_layout()
2468     plt.show()
2469
2470     # ----- Summary Table -----
2471     logger.info("\n--- Final Comparison Summary (Pandas) ---")
2472     summary = []
2473     for name, res in all_results.items():
2474         if isinstance(res, dict) and res.get("rmse"):
2475             best_rmse = min([v for v in res["rmse"] if np.isfinite(v)], default=np.nan)
2476             final_rmse = get_last_finite(res, "rmse")

```



```

2477         final_time = get_last_finite(res, "time")
2478         summary.append({"solver": name, "best RMSE": best_rmse, "final RMSE": final_rmse, "train time (s)": final_time})
2479     if summary:
2480         summary_df = pd.DataFrame(summary).sort_values("best RMSE")
2481         display(summary_df) # Use display for Colab/Jupyter
2482     else: logger.warning("No valid results to display in summary table.")
2483
2484     # ----- Save Figures -----
2485     FIG_DIR = Path(DATA_DIR_STR) / "figs"
2486     try:
2487         FIG_DIR.mkdir(exist_ok=True, parents=True)
2488         logger.info(f"Saving figures to {FIG_DIR}...")
2489         for n, fig_num in enumerate(plt.get_fignums(), 1):
2490             plt.figure(fig_num)
2491             fig_path = FIG_DIR / f"solver_plot_{n:02d}.png"
2492             plt.savefig(fig_path, dpi=180, bbox_inches='tight')
2493             print(" saved →", fig_path)
2494     except Exception as e_fig: logger.error(f"Could not save figures: {e_fig}")
2495
2496     print("\n✅ All comparisons finished.")
2497
2498 # --- ADDED Block 6-a: Run OT Demo (Rank 0 only) ---
2499 if RANK_MPI == 0 and OT_AVAILABLE:
2500     logger.info("\n+++ Running OT Barycentre Demo +++")
2501     try:
2502         if 'run_barycentre_demo' in globals():
2503             ot_demo_results = run_barycentre_demo()
2504             plt.figure(figsize=(6, 4))
2505             plt.plot(ot_demo_results['grid'], ot_demo_results['sources'], '--', label='Sources')
2506             plt.plot(ot_demo_results['grid'], ot_demo_results['barycenter'], 'r-', label='Barycenter')
2507             plt.title('Wasserstein Barycenter Demo'); plt.legend(); plt.tight_layout(); plt.show()
2508         else: logger.warning("run_barycentre_demo function not defined. Skipping OT demo.")
2509     except Exception as e_ot: logger.error(f"OT Barycentre Demo failed: {e_ot}")
2510
2511 # === ADDED Block 6: PCA Trajectory Plot (Rank 0 only) ===
2512 if RANK_MPI == 0 and PCA_AVAILABLE and TRAJECTORY_CACHE is not None and len(TRAJECTORY_CACHE) >= 3:
2513     logger.info("\n+++ Generating PCA Trajectory Plot +++")
2514     try:
2515         traj_dim = TRAJECTORY_CACHE[0].size
2516         valid_traj = [t for t in TRAJECTORY_CACHE if isinstance(t, np.ndarray) and t.size == traj_dim]
2517         if len(valid_traj) >= 3:
2518             pcs = PCA(n_components=2).fit_transform(np.vstack(valid_traj))
2519             plt.figure(figsize=(4.5,4)); plt.plot(pcs[:,0], pcs[:,1], '-o', markersize=3)
2520             plt.title('Optimisation Trajectory (PCA)'); plt.xlabel('PC1'); plt.ylabel('PC2')
2521             plt.tight_layout(); plt.show()
2522         else: logger.warning("Not enough valid trajectory points for PCA plot.")
2523     except Exception as e_pca: logger.error(f"PCA Trajectory plot failed: {e_pca}")
2524
2525

```

```

*** Requirement already satisfied: mpi4py in /usr/local/lib/python3.11/dist-packages (4.0.3)
Requirement already satisfied: POT in /usr/local/lib/python3.11/dist-packages (0.9.5)
Requirement already satisfied: numpy>=1.16 in /usr/local/lib/python3.11/dist-packages (from POT) (2.0.2)
Requirement already satisfied: scipy>=1.6 in /usr/local/lib/python3.11/dist-packages (from POT) (1.15.2)
Sun May 4 22:25:15 2025

```

NVIDIA-SMI 550.54.15				Driver Version: 550.54.15		CUDA Version: 12.4	
GPU	Name	Persistence-M	Bus-Id	Disp.A	Volatile	Uncorr. ECC	
Fan	Temp	Perf	Pwr:Usage/Cap	Memory-Usage	GPU-Util	Compute M.	MIG M.
0	NVIDIA A100-SXM4-40GB	Off	00000000:00:04.0	Off	0		
N/A	40C	P0	46W / 400W	0MiB / 40960MiB	0%	Default	Disabled

Processes:							
GPU	GI	CI	PID	Type	Process name	GPU Memory	
ID	ID					Usage	
No running processes found							

```

ERROR: Could not find a version that satisfies the requirement softimpute (from versions: none)
ERROR: No matching distribution found for softimpute
Mounted at /content/drive
+++ MPI Detected: Running with 1 processes. +++
+++ Mounting Google Drive +++
Mounted at /content/drive

```

Drive mounted.