# Primal Dual Hybrid Gradient Algorithm Enhancement Project

I'll help you brainstorm ideas and serve as your research assistant for your PDHG implementation. Let me break this down into actionable components for your two-week timeline.

## ROCmprofiler Analysis for Sparse vs Dense Matrices

### Setting Up ROCmprofiler for Algorithmic Analysis

ROCmprofiler offers several modes for analyzing your sparse vs dense PyTorch tensor implementations[1]:

**Basic Profiling Setup:**

```
# Profile with comprehensive metrics
rocprof-compute profile -n pdhg_analysis -- ./your_pdhg_script

# For detailed kernel analysis
rocprof --hip-trace --hsa-trace --timestamp on --basenames on ./your_script
```

**Key Metrics to Compare:**

- **Memory bandwidth utilization**: Critical for understanding sparse vs dense efficiency[2]
- **Kernel execution times**: Compare GEMM vs SpMV operations
- **Memory access patterns**: Sparse matrices often have irregular access patterns[3]
- **Cache efficiency**: Dense operations typically have better cache locality

### Specific Analysis Strategies

**1. Kernel-Level Comparison:**

```
# Filter specific kernels for comparison
rocprof-compute profile -k "gemm*" -k "spmv*" -- ./pdhg_script
```

**2. Memory Analysis:**

- Use `--list-metrics` to identify memory-related counters[1]
- Focus on L2 cache hit rates and memory controller utilization
- Compare memory throughput between sparse and dense operations

**3. Visualization Options:**

The ROCmprofiler output can be visualized using Chrome tracing or Perfetto[4], which will help you see timeline differences between sparse and dense operations.

## Understanding Gradient Landscape Geometry

### Dimensionality Reduction Techniques

For understanding the latent geometry of your gradient landscape, consider these approaches[5]:

**1. Principal Component Analysis (PCA):**

- Project high-dimensional gradient space onto 2-3D for visualization
- Identify principal directions of gradient variation
- Useful for understanding convergence patterns

**2. t-SNE/UMAP Projections:**

- Better for preserving local neighborhood structure
- Can reveal clustering patterns in gradient dynamics
- Particularly useful for understanding algorithm behavior across different problem instances

**3. Parallel Coordinates:**

- Represent each dimension as a vertical axis
- Show how gradient components evolve during iterations
- Effective for identifying correlated gradient components

### Advanced Geometric Analysis

**1. Hessian Eigenvalue Analysis:**

- Compute condition number evolution
- Identify ill-conditioned subspaces
- Guide preconditioning strategies

**2. Gradient Flow Visualization:**

- Create vector field plots in reduced dimensions
- Visualize convergence basins
- Identify saddle points and local minima

**3. Manifold Learning:**

- Apply techniques like Isomap or Locally Linear Embedding
- Understand intrinsic dimensionality of the solution space

- Reveal geometric structure not apparent in original coordinates

**Literature-Inspired PDHG Enhancements (2-Week Implementation)**

### Week 1: Core Algorithmic Improvements

**1. Adaptive Step Size Selection:**

- Implement Backtracking Line Search for primal/dual steps
- Add momentum-based step size adaptation
- Consider Nesterov-type acceleration schemes

**2. Preconditioning Strategies:**

- **For Sparse Matrices**: Implement diagonal preconditioning or incomplete factorizations
- **For Dense Matrices**: Use Cholesky or LU-based preconditioners
- Compare effectiveness using your ROCmprofiler analysis

**3. Convergence Acceleration:**

- Implement Anderson acceleration for the fixed-point iteration
- Add restart mechanisms based on gradient angle changes
- Consider quasi-Newton updates for the dual variables

### Week 2: GPU-Specific Optimizations

**1. Memory Access Optimization:**

- Implement tiled/blocked matrix operations for dense cases
- Use compressed sparse row (CSR) format optimizations for sparse matrices[3]
- Apply memory coalescing techniques based on your profiling results

**2. Mixed Precision Implementation:**

- Use FP16 for matrix operations where numerically stable
- Keep higher precision for accumulation and critical computations
- Leverage PyTorch's automatic mixed precision capabilities

**3. Multi-GPU Scaling:**

- Implement domain decomposition approaches
- Use asynchronous communication for gradient updates
- Consider parameter server architectures for large-scale problems

# Specific Research Directions

## Sparse Matrix Optimizations

Based on recent literature[6], focus on:

- **Load balancing**: Optimize work distribution across GPU cores

- **Binning strategies**: Group similar sparsity patterns for efficient processing

- **Numerical computation modes**: Optimize for your specific sparsity patterns

## Geometric Insights Implementation

### 1. Real-time Landscape Analysis:

```
# Pseudocode for gradient landscape tracking
def track_gradient_geometry(gradients, iteration):
    # Compute PCA of recent gradient history
    pca_projection = compute_pca(gradients[-window_size:])
    # Analyze curvature information
    hessian_approx = compute_bfgs_approximation(gradients)
    # Update visualization
    update_landscape_plot(pca_projection, hessian_approx)
```

### 2. Convergence Prediction:

- Monitor gradient alignment across iterations

- Track effective dimensionality of the gradient space

- Implement early stopping based on geometric criteria

## Implementation Priority Matrix

| Priority | Enhancement | Implementation Time | Expected Impact |
|----------|-------------|---------------------|-----------------|
| High | ROCmprofiler integration | 2-3 days | High - foundational for optimization |
| High | Adaptive step sizes | 2-3 days | High - improves convergence |
| Medium | Sparse matrix optimizations | 3-4 days | High - performance critical |
| Medium | PCA-based landscape analysis | 2-3 days | Medium - insights for tuning |
| Low | Mixed precision | 1-2 days | Medium - memory and speed gains |
| Low | Multi-GPU scaling | 4-5 days | High - but complex implementation |

### Recommended Tools and Libraries

**For Profiling and Analysis:**

- ROCmprofiler for detailed GPU analysis[1] [2]
- PyTorch's built-in profiler for integration
- Perfetto or Chrome tracing for visualization[4]

**For Geometric Analysis:**

- scikit-learn for PCA and manifold learning
- matplotlib/plotly for interactive visualizations
- PyTorch's autograd for Hessian approximations

**For Sparse Operations:**

- PyTorch sparse tensors with ROCm backend
- Consider rocSPARSE integration for critical operations[7]

This approach gives you a solid foundation for both understanding your algorithm's behavior through profiling and implementing meaningful enhancements within your two-week timeframe. Focus on the high-priority items first, using ROCmprofiler analysis to guide your optimization decisions.

❄

1. https://rocm.docs.amd.com/projects/rocprofiler-compute/en/latest/how-to/use.html
2. https://rocm.docs.amd.com/projects/rocprofiler/en/latest/how-to/using-rocprof.html
3. https://gpuopen.com/learn/amd-lab-notes/amd-lab-notes-spmv-docs-spmv_part1/
4. https://extremecomputingtraining.anl.gov/wp-content/uploads/sites/96/2022/11/ATPESC-2022-Track-6-Talk-4-Tandon-AMD.pdf
5. https://www.reddit.com/r/learnmath/comments/1fki7yh/how_can_i_visualize_the_solutions_of_gradient/
6. https://link.springer.com/article/10.1007/s11227-024-06234-2
7. https://rocm.docs.amd.com/projects/rocSPARSE/en/develop/reference/enumerations.html