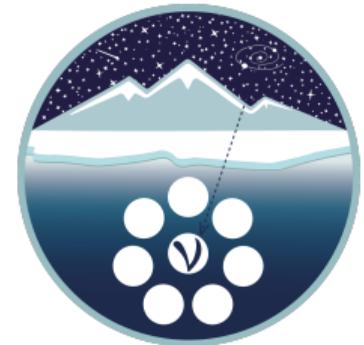


# Multi-Threading for Baikal-GVD Core Software Framework

A.G. Solovjev  
on behalf of the Baikal-GVD Collaboration

Joint Institute for Nuclear Research, Dubna, Russia, 141980

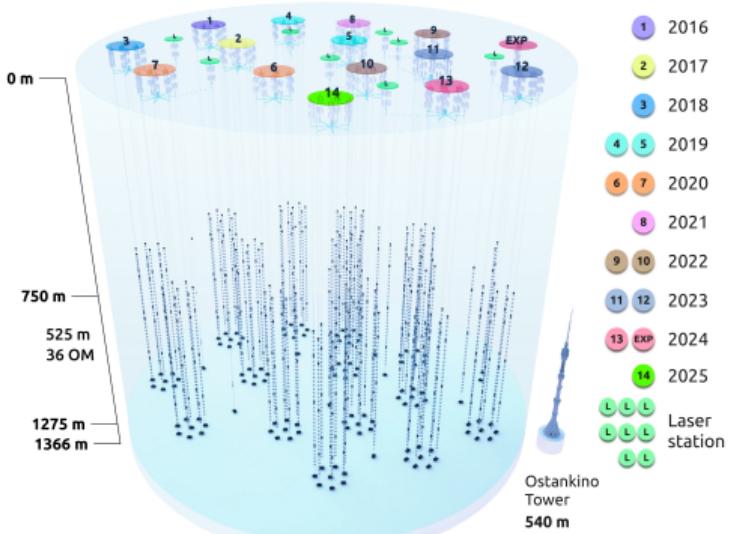
Distributed Computing and Grid-technologies  
in Science and Education, 7-11 July 2025





# Introduction

## Baikal-GVD neutrino telescope in Lake Baikal



14 clusters ( $0.7 \text{ km}^3$ ), 8 strings/cluster, 36 OMs/string.  
Expanding yearly (+1–2 clusters/year).

## Data Processing System

- **Core:** BARS framework (C++)
- **Management:** Python workflow

## Parallelism Levels

- Cluster-level (distributed VMs)
- Workflow-level (parallel processing sections):
  - Fast (2–13 min/file)
  - Offline (1–5 hr/run)
- **Thread-level (this work)**



A. G. Solovjev et al. (Baikal-GVD), "Review of Software for Automatic Processing and Analysis of Data from the Baikal-GVD Neutrino Observatory". In: *Parallel Computational Technologies*. Ed. by L. Sokolinsky et al. Springer, 2024, pp. 80–91. [https://doi.org/10.1007/978-3-031-73372-7\\_6](https://doi.org/10.1007/978-3-031-73372-7_6)

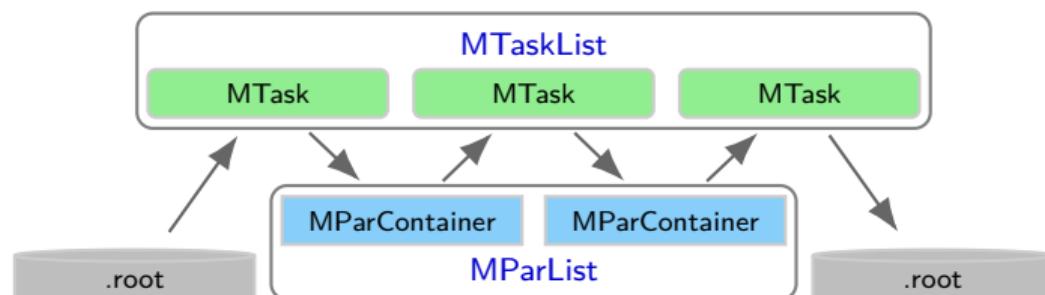


## Processing Workflow

### 1. Pre-processing:

- Register containers
- Load parameters

### 2. Processing:



### 3. Post-processing:

- Finalize results
- Release resources

## Event Loop Code

```
Int_t rc = kTRUE;
// Phase 1: Pre-process all tasks
for (MTask *task : fTaskList) {
    rc = task->PreProcess(fParList);
}
// Phase 2: Main processing loop
if (rc) {
    Parallelization Target
    while (true) {
        fParList->Reset();
        for (MTask *task : fTaskList) {
            rc = task->Process();
        }
        if (!rc) break;
    }
}
// Phase 3: Post-process all tasks
for (MTask *task : fTaskList) {
    rc = task->PostProcess(fParList);
}
return rc;
```

# Parallelization Requirements and Performance



## Key Constraints

- **Event Processing Order** must be strictly preserved
- **Task Execution Order** within tasklist must be maintained
- **Task Exclusivity**: identical tasks must not run simultaneously

## Theoretical Speedup

$$T_{\text{sequential}} = \bar{t}_{\text{tasklist}} \cdot N_{\text{events}}$$

$$T_{\text{parallel}} \approx (\bar{t}_{\text{tasklist}} + \bar{t}_{\text{idle}}) \cdot \frac{N_{\text{events}}}{n_{\text{threads}}}$$

$$S \approx \frac{\bar{t}_{\text{tasklist}}}{\bar{t}_{\text{tasklist}} + \bar{t}_{\text{idle}}} \cdot n_{\text{threads}}$$

## Amdahl's law

$$S_{\text{total}} = \frac{1}{(1-p) + \frac{p}{S}}, \quad p = \frac{T_{\text{Process}}}{T_{\text{total}}}$$

## OpenMP Tasking

```
// Thread pool setup
#pragma omp parallel
#pragma omp single
while (true) {
    #pragma omp task
    {
        // Process tasklist
        fParList->Reset();
        for (MTask *task : fTaskList) {
            omp_set_lock(task->fLock);
            rc = task->Process();
            omp_unset_lock(task->fLock);
        }
        if (!rc) break;
    }
}
```



# Optimization Strategies

## Core Approaches to Reduce Idle Time ( $\bar{t}_{idle}$ )

- **Task Decomposition**

- Split long task into smaller subtasks
- Advantages:
  - ▶ Enables dynamic load balancing
  - ▶ Minimizes synchronization overhead

- **Task Optimization**

- Optimize individual task execution time
- Methods:
  - ▶ Profile and optimize critical paths
  - ▶ Implement strategic caching

## Why Many Small Tasks Work Better

- Few large tasks (causes thread starvation)



- Many small tasks (maximizes core utilization)



- Advantages

- More tasks  $\Rightarrow$  Better load balancing
- Smaller tasks  $\Rightarrow$  Less idle time

## Approaches to Increase Parallel Portion ( $p$ )

- Pre/Post-Process Optimization

- Offload critical paths to separate processes
- Pre-cache frequently used data



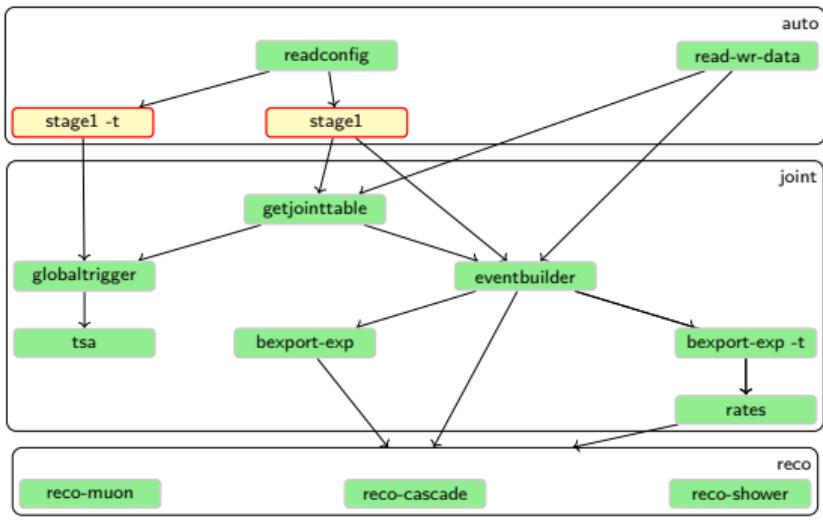
# Optimized Programs by Processing Stage

Basic parallelization does not speed up enough

## Per-File Workflow

Key optimized program:

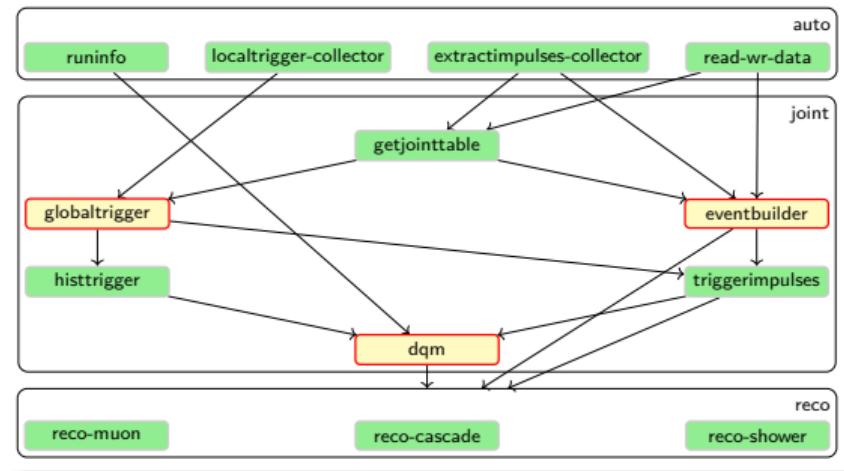
- stage1



## Per-Run Workflow

Key optimized program:

- globaltrigger/eventbuilder
- dqm

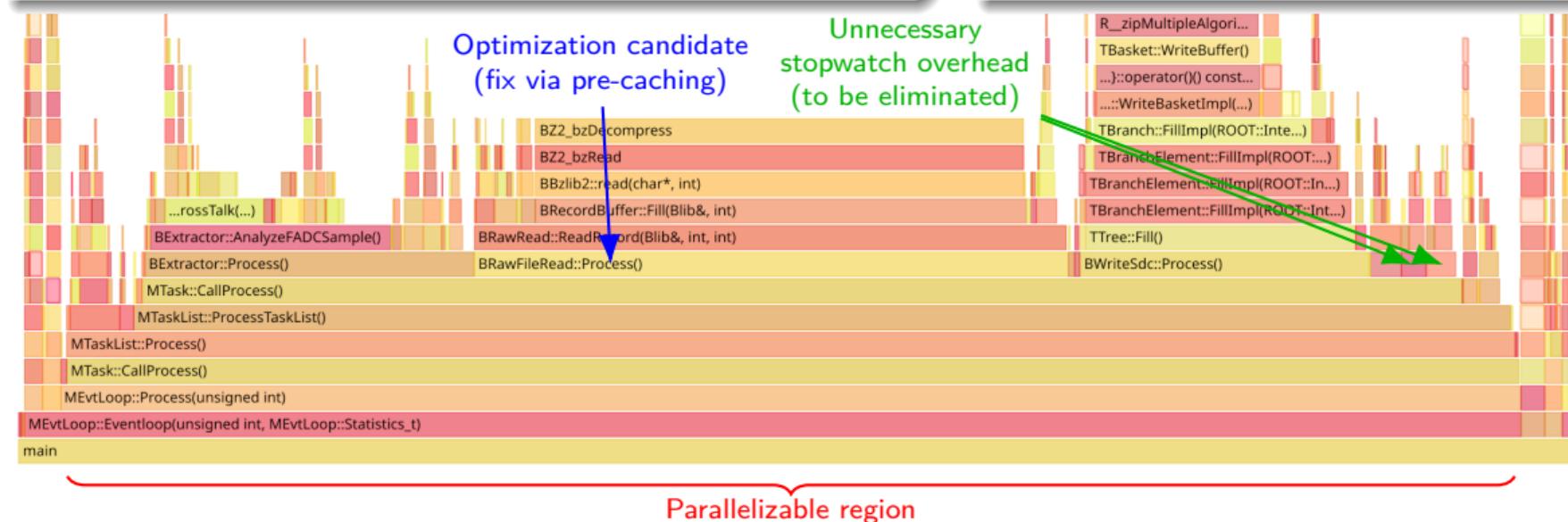




# Flame Graph Analysis for stage1

## Profiled Command

```
stage1  
--season 2024 --cluster 2 --run 230 --file 1  
--create-chain --absolute-time --raw-write-additional  
(261124 events)
```



## Flame graphs

**Intuitive visualization** of profiling:

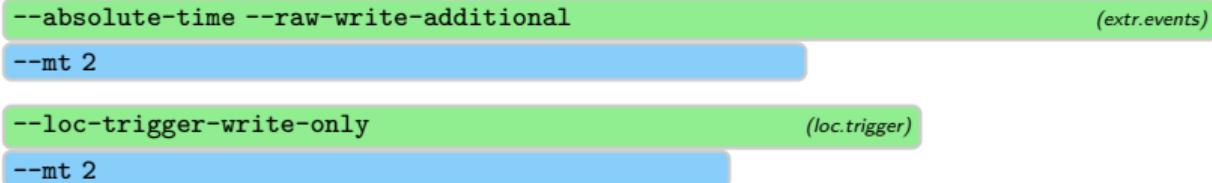
- element width — execution time
- hierarchy — call stacks



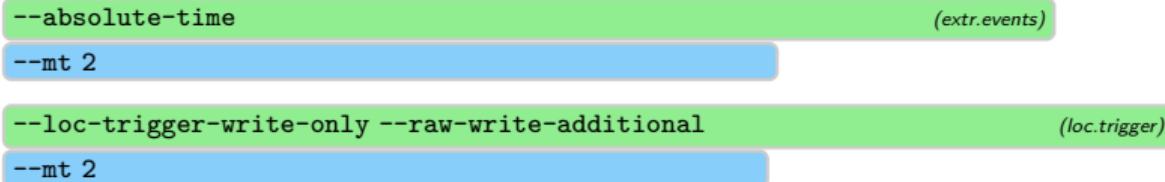
# Parallelization Strategies for stage1

```
stage1 --season 2024 --cluster 2 --run 230 --file 1 --create-chain ...
```

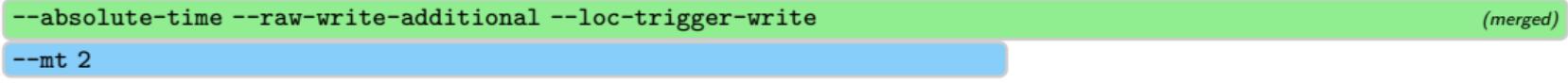
- **Original parallel execution of independent instances (MT faster, 4 cores)**



- **Parallel execution with optimized options combinations (MT faster, 4 cores)**



- **Merged execution (MT no slowdown, 2 cores)**



0

10s

20s

30s



# Parallelization Bottlenecks in globaltrigger/eventbuilder



## Execution Details

- **Command:**

```
globaltrigger  
--season 2024 --cluster 2 --run 230  
--mt 2 [2 threads]
```

- **Execution Log:**

Main Tasks:

```
88.9% BGlobTriggerBuilder (idle 47.4%)  
9.0% BWriteTree (idle 36.3%)
```

## Performance

- Efficiency:  $\approx 54\%$
- Speedup: —

## Solution

- Decompose monolithic SDC processing into per-SDC tasks

## Current Implementation

```
class BGlobTriggerBuilder : public MTask {  
    Int_t Process(int onThread) {  
        for (auto sdc : SDCs) {  
            // Processes ALL SDCs sequentially  
        }  
    }  
};
```

## Core Problem

- **Sequential** SDC processing in a loop
- **Self-blocking** in multi-threaded execution



# Task Decomposition in globaltrigger/eventbuilder

## Improved Implementation

```
// Per-SDC parallel tasks
class BGlobTriggerBuilderSdc : public MTask {
    Int_t Process(int onThread) {
        // Processes individual SDC
    }
};

// Main container
class BGlobTriggerBuilder : public MTaskList {
    BGlobTriggerBuilder() {
        for (auto sdc : SDCs) {
            AddToList(new BJointEventBuilderSdc(sdc));
        }
    }
};

// User code
- tasks.AddToList(&globtrigger);
+ tasks.AddToList(*globtrigger.GetList());
```

## Execution Log

### SDC Tasks:

```
2.5% BGlobTriggerBuilder190 (idle 1.3%)
2.0% BGlobTriggerBuilder192 (idle 1.7%)
1.9% BGlobTriggerBuilder193 (idle 1.5%)
...
3.2% BGlobTriggerBuilder215 (idle 1.3%)
```

## Parallel Performance

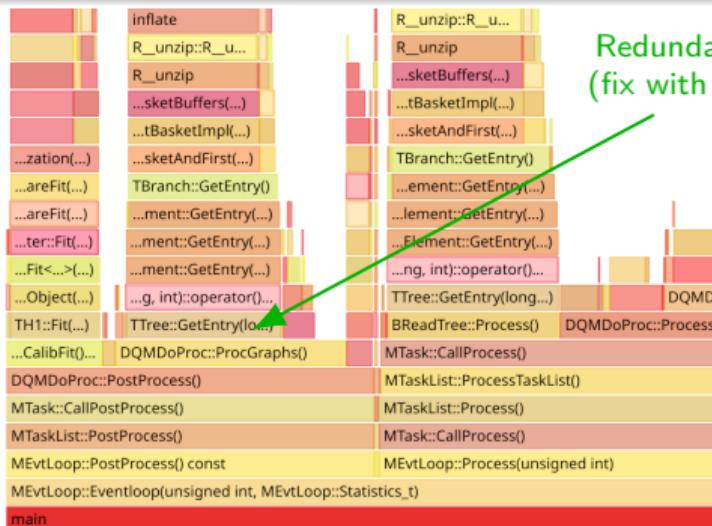
	2 cores	3 cores	4 cores
• Efficiency:	<b>92%</b>	<b>86%</b>	<b>81%</b>
• Speedup:	<b>1.7×</b>	<b>2.3×</b>	<b>2.8×</b>



# Flame Graph Analysis for DQM

## Profiled Command

```
dqm-cli
--season 2024 --cluster 9 --run 112
--exponent --poisson --uniformity --tdqm --charge
```



Redundant tree pass  
(fix with pre-caching)

Candidate for  
task extraction

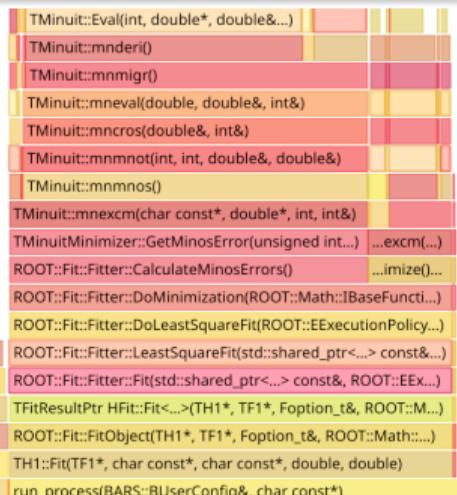
Parallelizable region

External process candidate

## Execution Time

≈ 4 minutes

(72986 events)



External process candidate



# Optimization Results for DQM

## Profiled Command

```
dqm-cli  
--season 2024 --cluster 9 --run 112  
--exponent --poisson --uniformity --tdqm --charge  
--mt 2
```

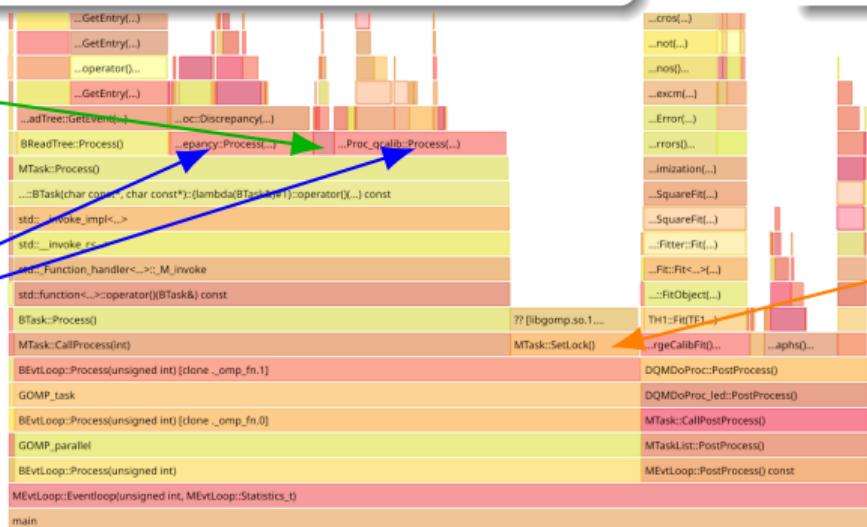
[2 threads]

## Speedup

≈ 2× faster  
(4 min ↓ 2 min)

Pre-caching  
tree pass

Long task  
decomposition



Remaining  
synchronization  
overhead

Parallelized region (efficiency ≈75%)



# Conclusion

## Key Achievements

- Successfully implemented multi-threading in Baikal-GVD's core framework:
  - **Eventloop parallelization** via OpenMP tasking
  - **Optimization strategies** via task granularity
- Achieved significant speedups across processing stages:
  - stage1: stable performance during combined execution without increased resource consumption
  - globaltrigger/eventbuilder: good scaling via parallel detector section processing
  - dqm: 2× speedup via improved load balancing and elimination of sequential bottlenecks

## Future Directions

- **Refactoring processing graphs** to enable multithreading support
- **Eliminating exclusivity requirement** through tasks refactoring



# Conclusion

## Key Achievements

- Successfully implemented multi-threading in Baikal-GVD's core framework:
  - **Eventloop parallelization** via OpenMP tasking
  - **Optimization strategies** via task granularity
- Achieved significant speedups across processing stages:
  - stage1: stable performance during combined execution without increased resource consumption
  - globaltrigger/eventbuilder: good scaling via parallel detector section processing
  - dqm: 2× speedup via improved load balancing and elimination of sequential bottlenecks

## Future Directions

- **Refactoring processing graphs** to enable multithreading support
- **Eliminating exclusivity requirement** through tasks refactoring

**Thank you!**