Saint Petersburg State University
Department of computer modeling and multiprocessor systems
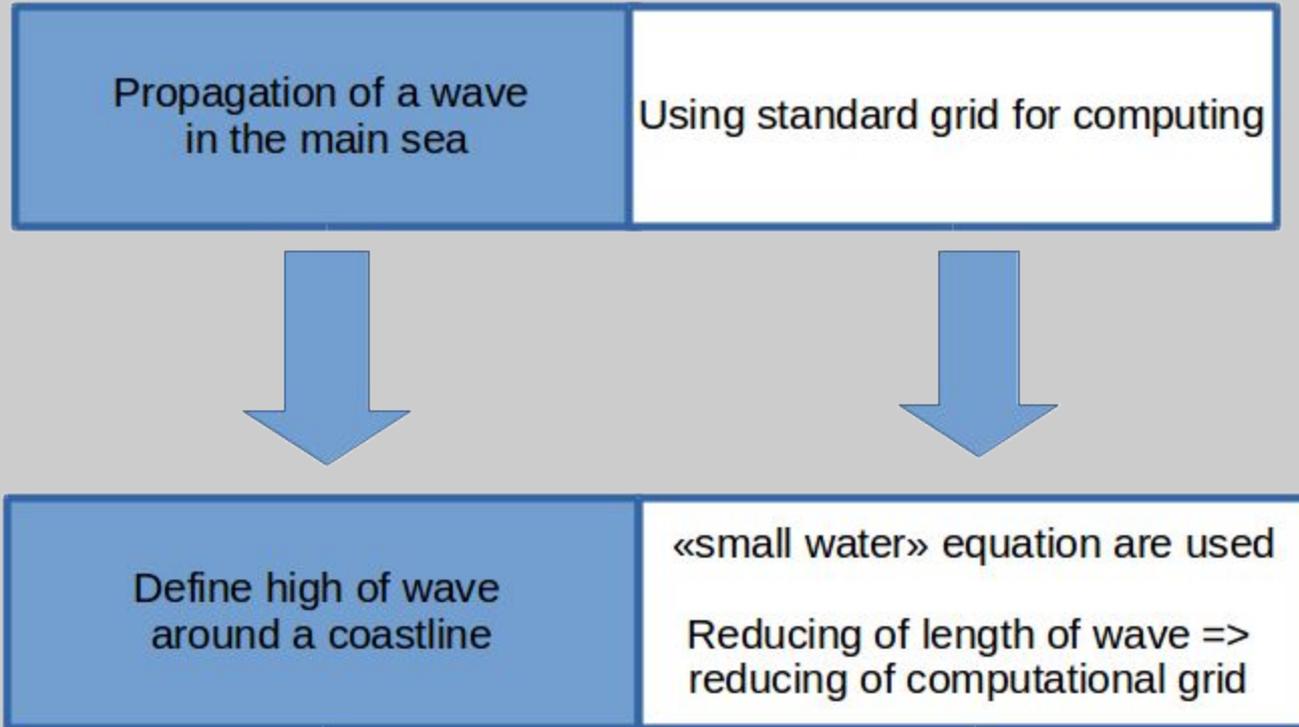
**Svetlana Sveshnikova**

**Processing of multidimensional data in distributed systems for solving the task of tsunami waves modeling**
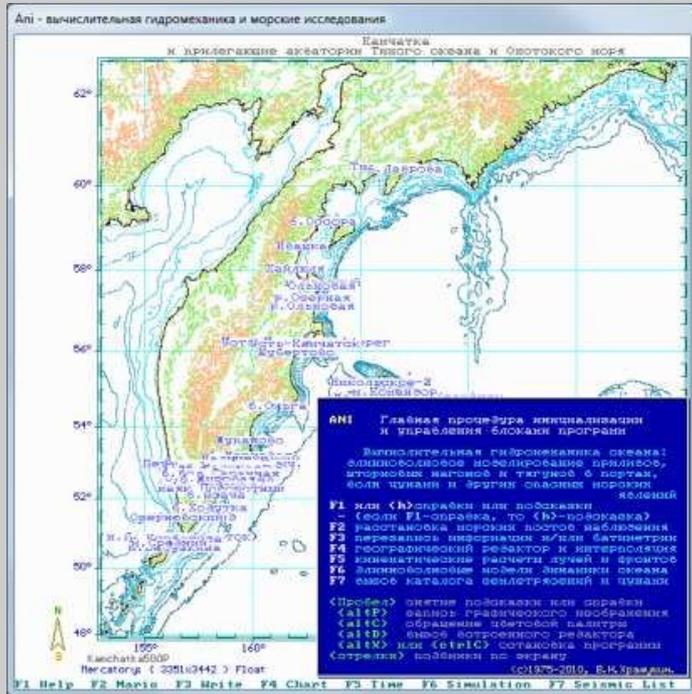
GRID' 2016

# Introduction

There are 2 stages of waves propagation and 2 computer modeling methods.

| | |
|---|---|
| Propagation of a wave in the main sea | Using standard grid for computing |
| Define high of wave around a coastline | «small water» equation are used<br><br>Reducing of length of wave => reducing of computational grid |

# Existing solution



- Works only for one node
- Bathymetry was full downloaded in RAM
- Only for Windows

Храмушин В. Н. Прямые вычислительные эксперименты для моделирования цунами, штормовых нагонов, экстремальных течений и приливного режима в открытом океане и вблизи побережья (г/р № 2010615848).

# Our task

To develop tools for operative processing grid re-interpolation indicated area from bathymetry files for solve modeling tsunami tasks.

**System requirements:**

1. NetCDF files must have processing
2. Work with 14 gb files in operative memory
3. Selection of given section on the map for further operations (coordinates and required accuracy given by user)

# Why distributed computing?

✓ High speed
✓ Availability for big data processing
✓ Reliability and fault-tolerance

Our choice it is framework Apache Spark

Apache Spark implements that structure as RDD.

RDD - resilient distributed dataset for speed big data processing in operating memory.



| Transformations | Actions |
| --- | --- |
| map (func) | reduce(func) |
| flatMap(func) | collect() |
| filter(func) | count() |
| groupByKey() | first() |
| reduceByKey(func) | take(n) |
| mapValues(func) | saveAsTextFile(path) |
| sample(…) | countByKey() |
| union(other) | foreach(func) |
| distinct() | … |
| sortByKey() | |
| … | |

# Some problems...

✓ Spark is designed for streaming and non-structure data processing

✓ NetCDF format use metadata and multidimensional arrays and require random access to file

**NetCDF file**

**metadata**

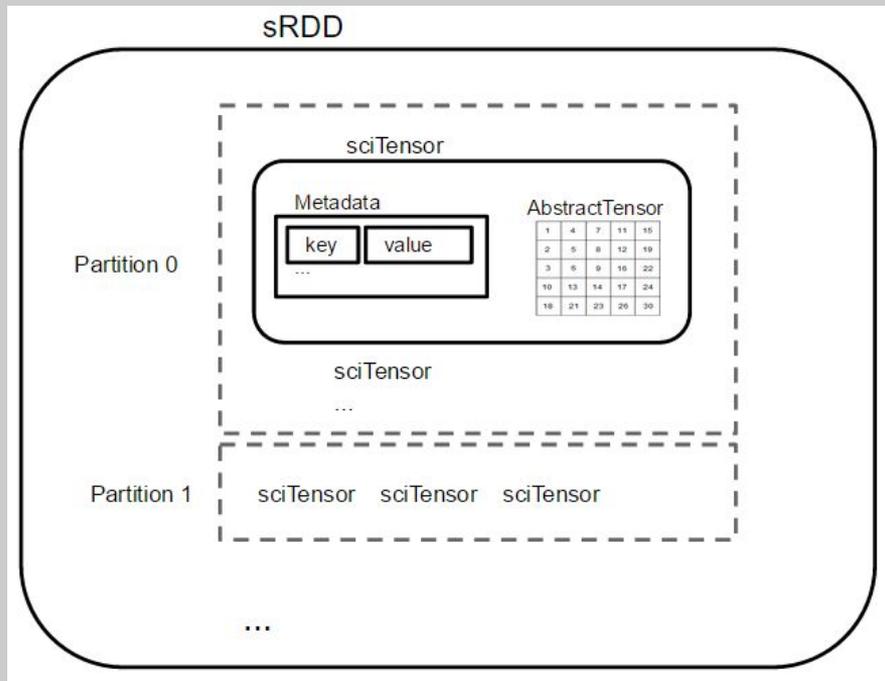Dimensions: ...
Variables: ...
Global atributes: ...

**data**

X = 3, -4, 0, 23, -12, 5, … 13;
Y = 24, 75, -23, 90, … 45;
…

# One of solutions

At the moment there are several solutions for processing data in NetCDF format on big-data systems (Hadoop and Spark). For Spark it is SciSpark library.

SciSpark - project supported by Apache Foundation and NASA Laboratory, that works with NetCDF files used linear algebra libraries (Breeze and ND4J).

sRDD - distributed dataset, oriented on the scientific data processing, in particular NetCDF.

# Advantages and disadvantages

- ☻ First available solution for processing geodata in Spark
- ☻ Support metadata and multidimensional arrays

- ☹ No official documentation: mans, tutorials, etc.
- ☹ Raw state of the product

# Data
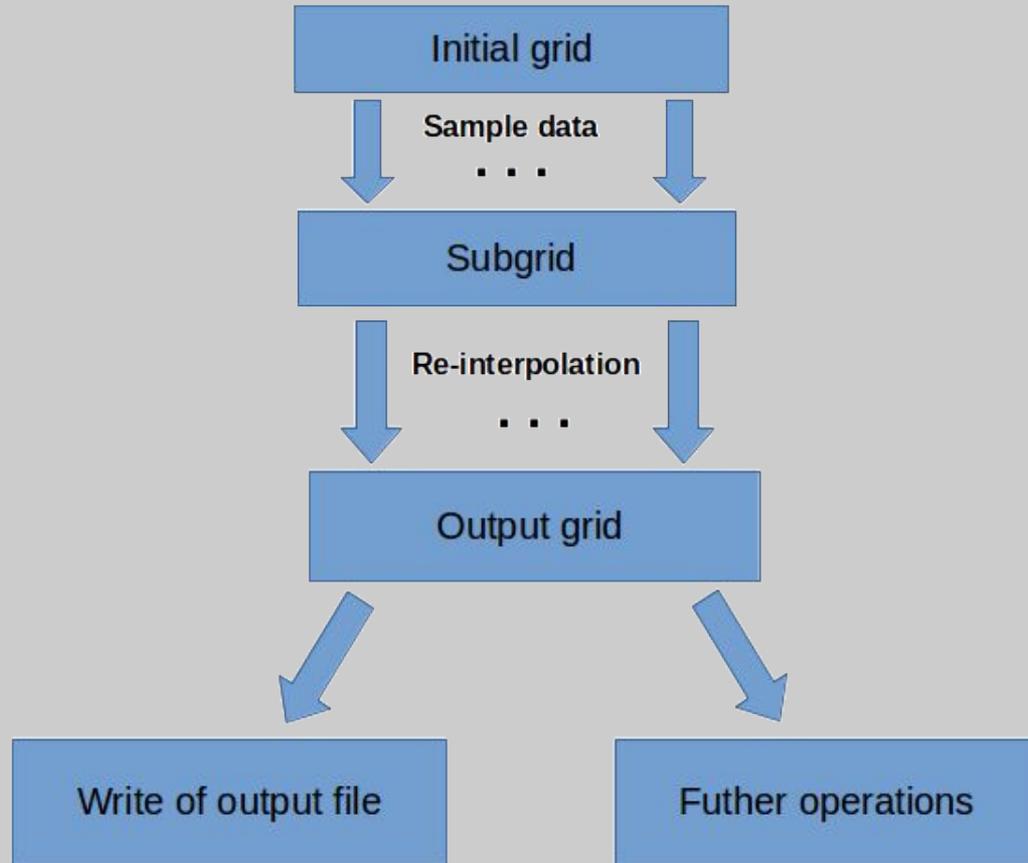
Data for processing: file with Earth's topography

Source: [ftp://topex.ucsd.edu/pub/srtm15_plus/](ftp://topex.ucsd.edu/pub/srtm15_plus/)

Data Size: 14 Gb

Step of the grid: 15 seconds

Data origin: the data given by processing of result altimetry survey from satellites CryoSat-2, Jason-1 and etc.
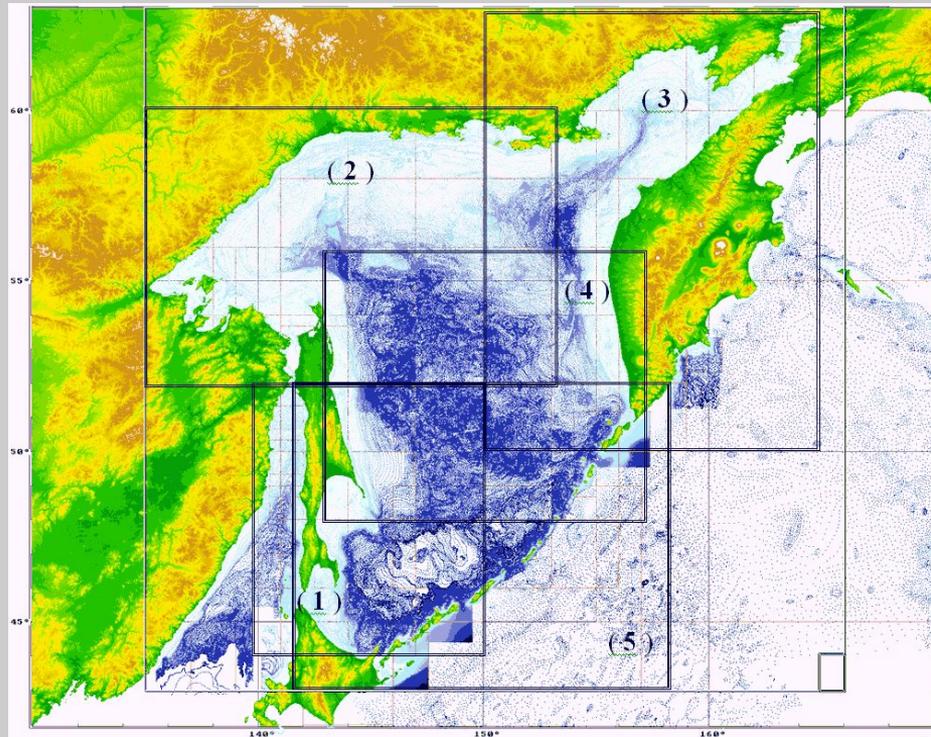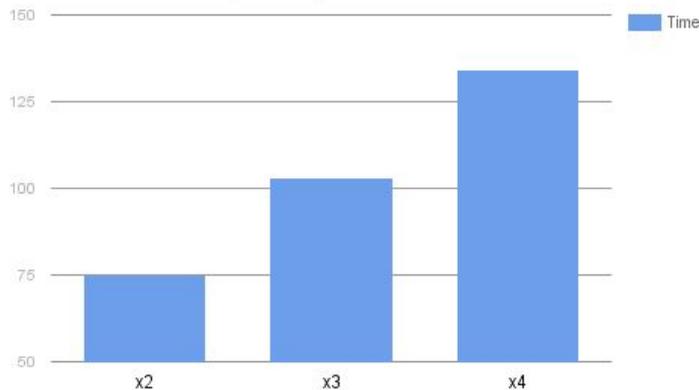
# Block-scheme of the program

# Cluster configuration

| CPU | Intel Xeon E5440, 2,83GHz |
|---|---|
| RAM | 4Gb |
| HDD | ST3250310NS, 7200 |
| Number of nodes | 12 |
| Number of cores per node | 8 |
| Software | Spark 1.6.0 + GlusterFS 3.6.3 |

# Results

1 file ~50Mb. Grid 30x15 degrees.
(7200x3600). Sample for interpolation
10x10 degrees.

# Thank you for your attention!
## Questions?