# Data-intensive Analysis and High Performance Computing in Solid Earth Sciences

Integrated Data Infrastructures

Networks

Computational Infrastructures

Data-Intensive Science

Jean-Pierre Vilotte, IPGP (CNRS-INSU)

With the contributions of: Dimitri Komatitsch (LMA, Marseille), Jean Virieux (ISTerre, Grenoble), N. Shapiro (IPG Paris), Eléonore Stutzmann (IPG Paris), Alexandre Fournier (IPG Paris), and the VERCE Team

**VERCE**

**EPOS** EUROPEANPLATEOBSERVINGSYSTEM

Beijing, May 16-21, 2013

# Data-intensive Research

## International community

- Global observation and monitoring systems
- Integrated Distributed Data Archives
- Data and metadata format standards
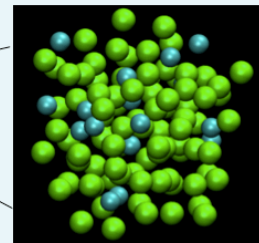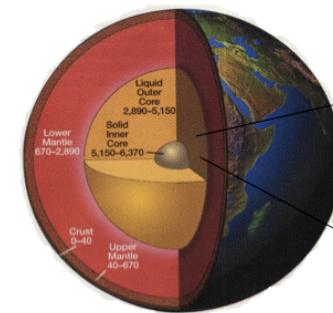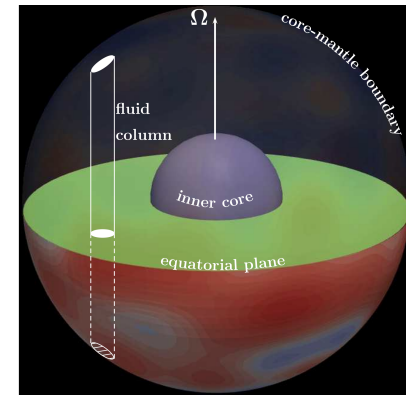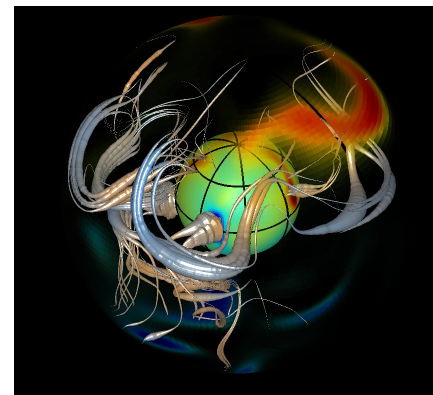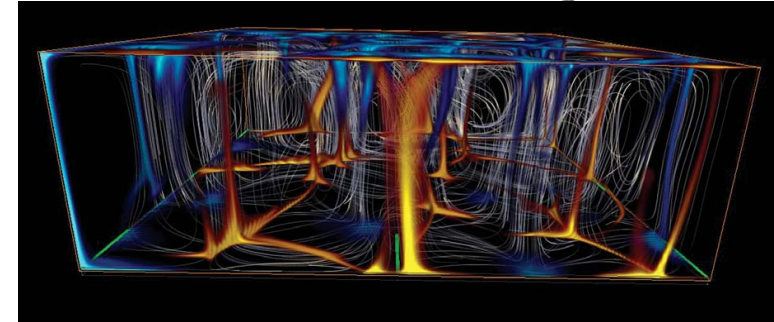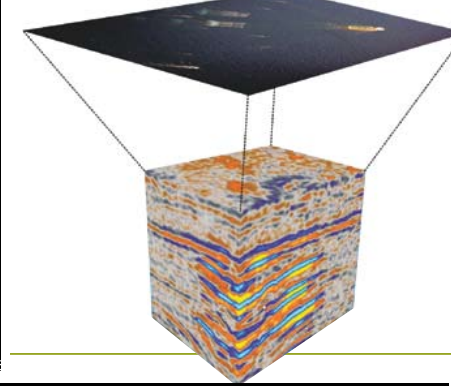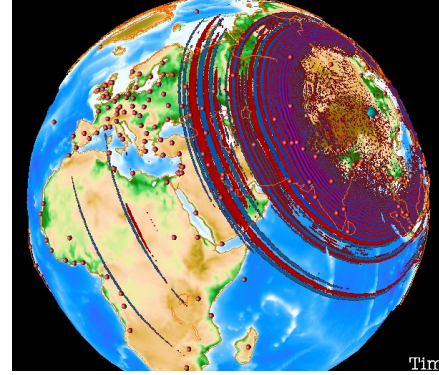
## Scientific challenges

- Understanding Earth's dynamics and structures
- Imaging Earth's interior and seismic sources

## *Augmented societal applications*

- Natural hazard and risk mitigation;
- Energy resources exploration and exploitation;
- Underground wastes and carbon sequestration;
- Nuclear test monitoring and treaty verification

## Data-intensive computing challenges

- Source detection and waveform data analysis
- High resolution inversion and data assimilation
- Quantification of forward/inverse uncertainties

# Computational Chalenges

**Massive data sets generated from observation systems and numerical simulations**
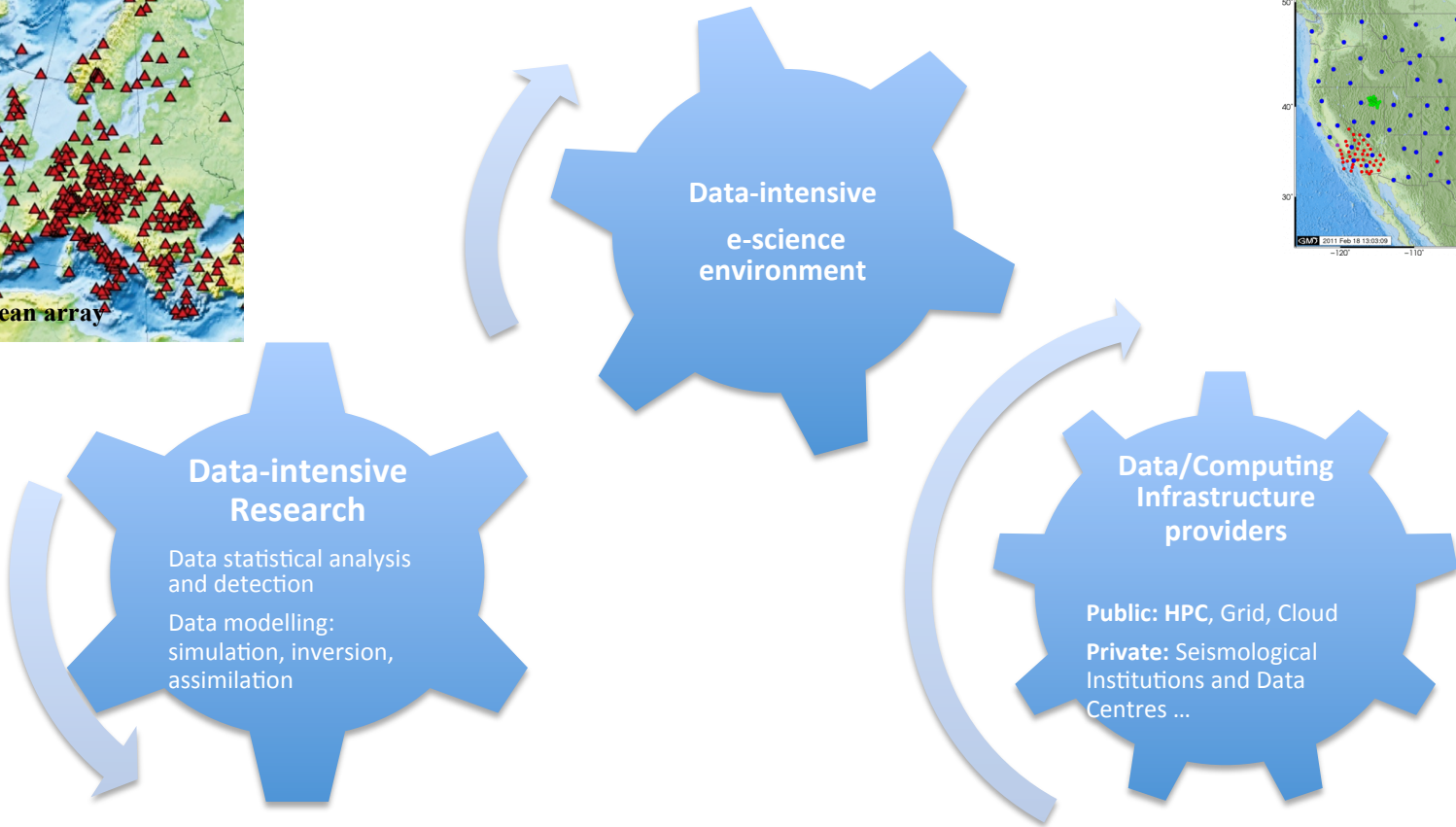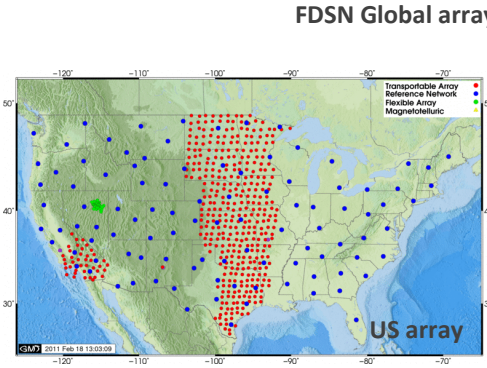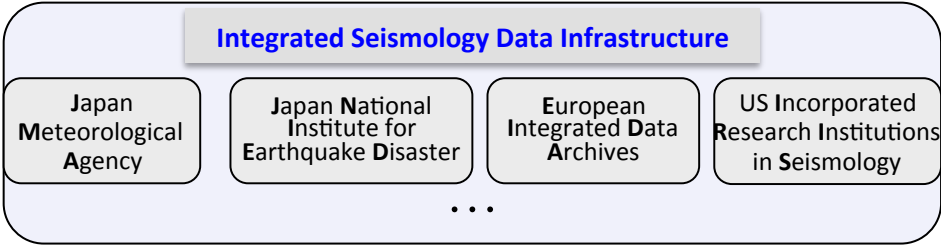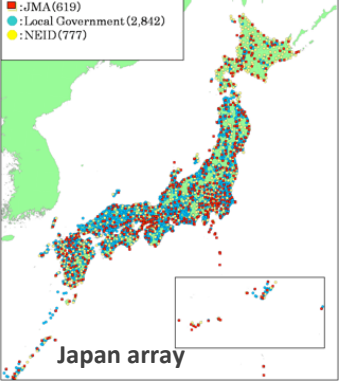
① **Data intensive statistical analysis:**

- **Monitoring property variations**: seismic noise correlation …
- **Seismic sources detections**: Coherent Interferometry (seismic and geodesy) …
- **Complex data processing**: GPS analysis, InSAR, optical image correlation analysis …

② **Data intensive modeling applications:**

- **Inversion (adjoint methods)**: geodynamo, acoustic and seismic full waveform tomography
- **Quantifying inverse uncertainties** (Monte Carlo): Tomography, geodesy, earthquake imaging
- **Time lapse tomography**: exploration seismology
- **Coherent interferometry and noise correlation tomography**: seismic tomography/migration, time reversal, seismic source imaging
- **Data assimilation**: geodynamo, seismic source imaging, mantle convection

③ **CPU intensive applications:**

- **Multi-physics simulations**: core-mantle dynamics, geological climate evolution, acoustic/elastodynamics coupling, tsunami/seismic sources
- **Multi-scales simulations (homogenization)**: wave propagation, earthquake dynamics, geodynamo
- **Stochastic quantification of forward uncertainties and variability**: geological climate evolution, wave propagation, earthquake dynamics, geodynamo
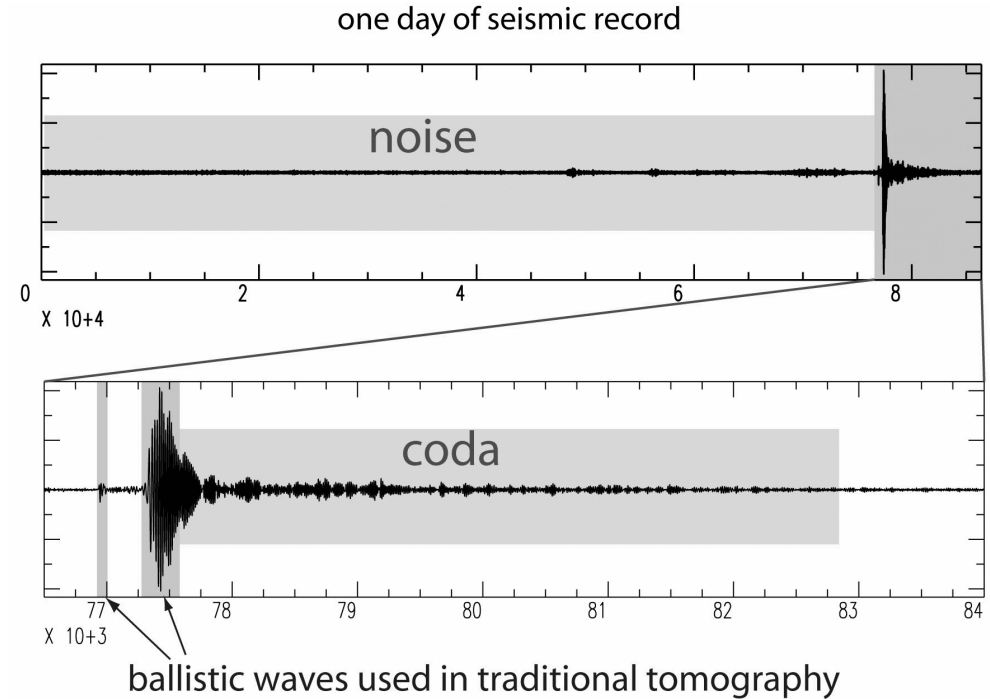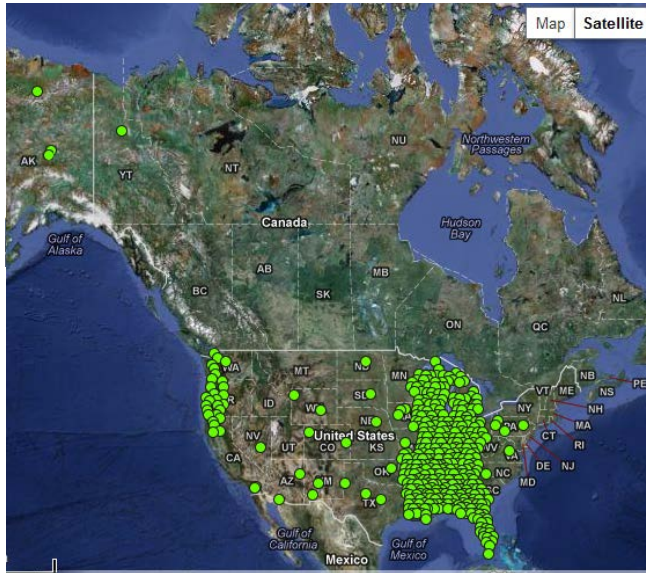
**Japan array**

**European array**

**Integrated Seismology Data Infrastructure**

| **J**apan **M**eteorological **A**gency | **J**apan **N**ational **I**nstitute for **E**arthquake **D**isaster | **E**uropean **I**ntegrated **D**ata **A**rchives | US **I**ncorporated **R**esearch **I**nstitutions in **S**eismology |
| --- | --- | --- | --- |

· · ·

**FDSN Global array**

**US array**

**Data-intensive e-science environment**

**Data-intensive Research**

Data statistical analysis and detection

Data modelling: simulation, inversion, assimilation

**Data/Computing Infrastructure providers**

**Public: HPC**, Grid, Cloud

**Private:** Seismological Institutions and Data Centres ...

egi

PRACE

*Cloud*

EP S EUROPEAN PLATE OBSERVING SYSTEM

| **Earth's interior imaging and dynamics: noise correlation, waveform analysis** | **Natural hazards: new tools for monitoring earthquakes, volcanoes, and tsunami** | **Interaction of solid Earth with Ocean and Atmosphere: environment, climate changes** |

# Data-Intensive statistical analysis: Seismic noise correlation



one day of seismic record

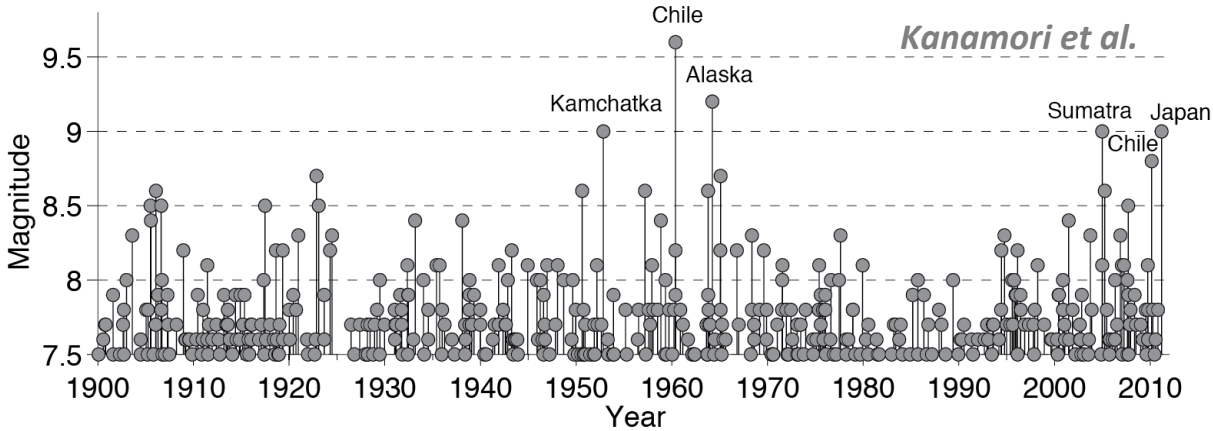noise

coda

ballistic waves used in traditional tomography

**Exploiting the statistical coherence in space and time of continuous waveforms records from dense arrays of broadband and strong motion instruments**
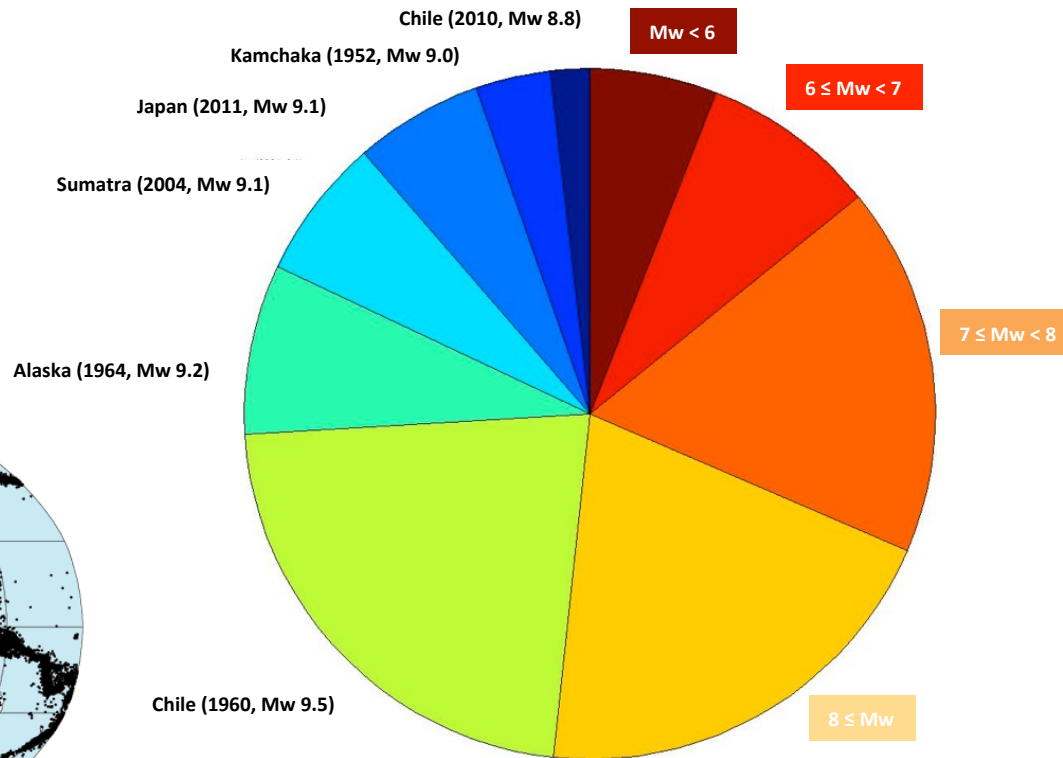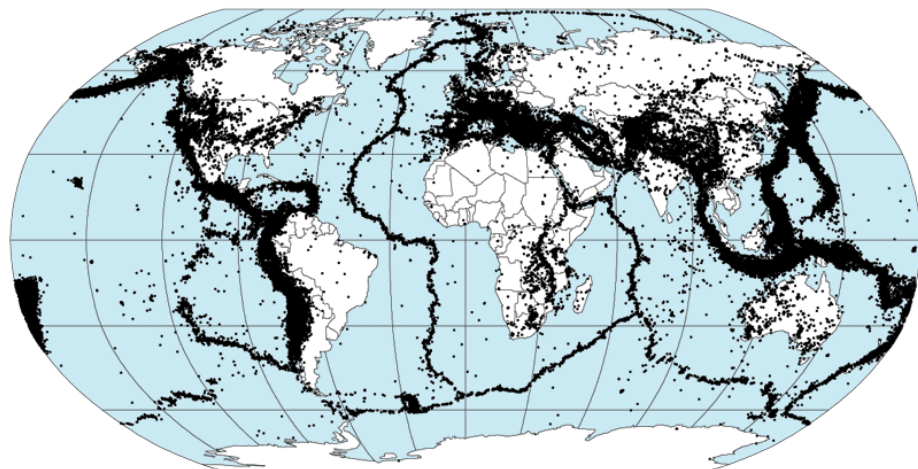
$$D = S \otimes M$$

**D** - seismic data
**S** - seismic source
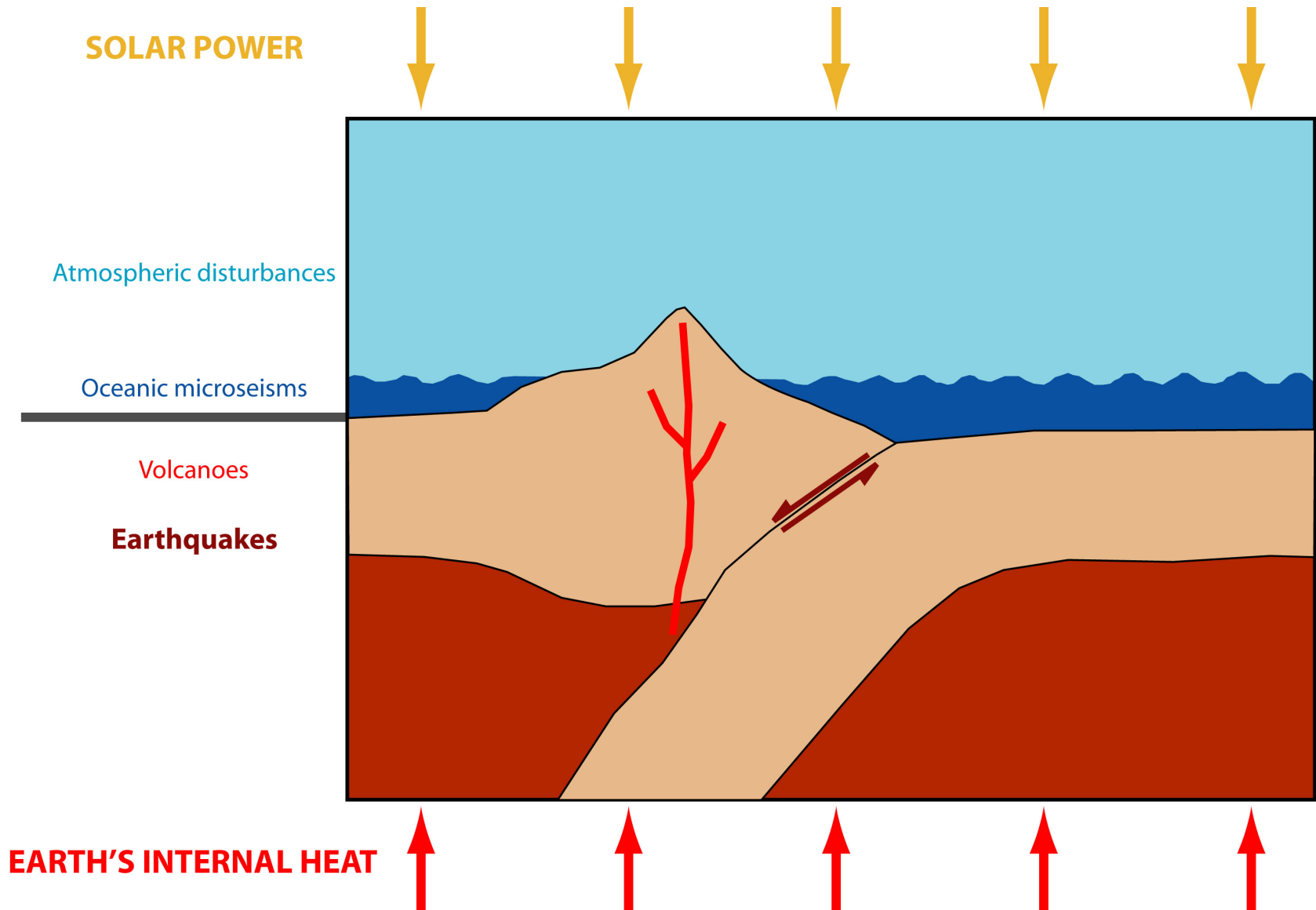**M** - media (Earth)

# Classical seismic sources: earthquakes



*Kanamori et al.*

Chile
Alaska
Kamchatka
Sumatra  Japan
Chile

Chile (2010, Mw 8.8)
Kamchaka (1952, Mw 9.0)
Japan (2011, Mw 9.1)
Sumatra (2004, Mw 9.1)
Alaska (1964, Mw 9.2)
Chile (1960, Mw 9.5)

Mw < 6
6 ≤ Mw < 7
7 ≤ Mw < 8
8 ≤ Mw

Preliminary Determination of Epicenters
358,214 Events, 1963 - 1998

# A wide range of natural seismic sources



SOLAR POWER

Atmospheric disturbances

Oceanic microseisms

Volcanoes

Earthquakes

EARTH'S INTERNAL HEAT

# Extracting Green function from random wave fields

For a **random** wave field with **homogeneous sources distribution** *everywhere* in the medium, it can been shown that:

$$\frac{d}{d\tau} C_{A,B}(\tau) = \frac{-\sigma^2}{4\,a} \left( G_a(\tau, \vec{r}_A, \vec{r}_B) - G_a(-\tau, \vec{r}_A, \vec{r}_B) \right)$$

**noise cross-correlation**                    **Green function**

- ✓ computing cross-correlation of seismic noise between two stations from long enough records is equivalent to an experiment when a source is acting at location of one of stations and recorded at another

- ✓ repetitive computations of noise cross-correlations are equivalent to using repetitive seismic sources and can be used to detect changes in the medium

$$D = S \otimes M \quad \Longrightarrow \quad$$

$$C(D,t) \approx M(t)$$

$$C(D,t) = Sc(t) \otimes M(t)$$

Helioseismology: Duvall et al. (1993)….; Laboratory Acoustics: Weaver and Lobkis (2001)…; Seismic coda waves: Campillo and Paul (2003)…; Marine acoustics: Roux et al., (2003)…; Ambient seismic noise: Shapiro and Campillo (2004)…
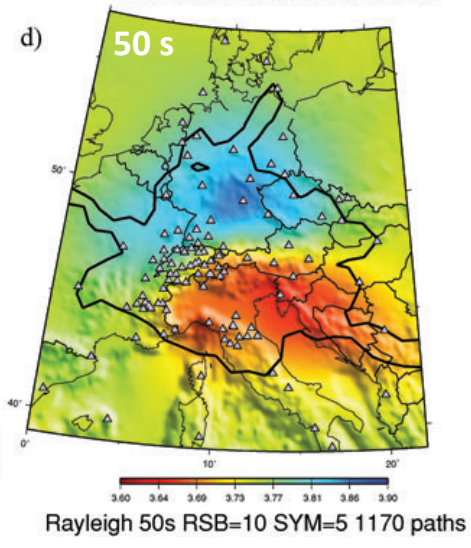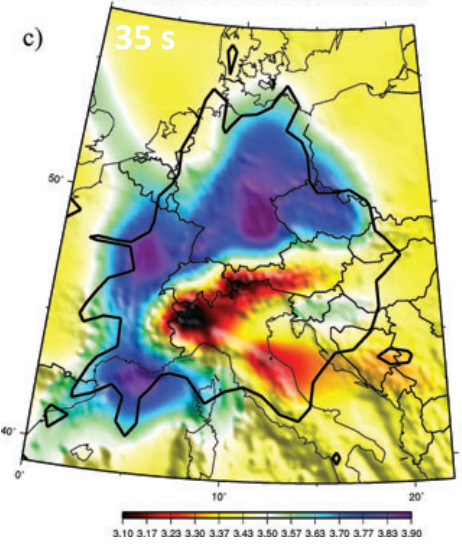
# Dense networks: local tomography

15-30s band-passed



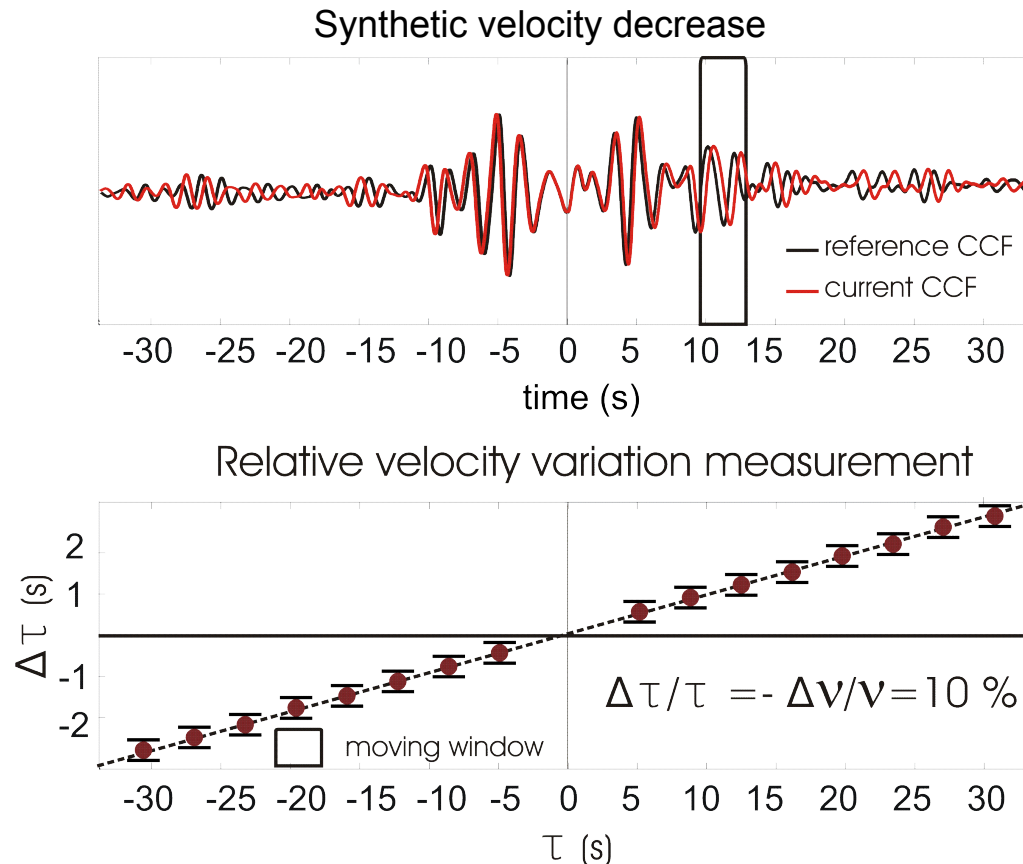Eikonal equation

$$\frac{\hat{k}_i}{c_i(r)} = \nabla\tau(r_i, r).$$

24 s Rayleigh wave

Stehly et al. (2009)

In the case of a homogeneous velocity perturbation in the media, waves travel times change proportionally to this perturbation.

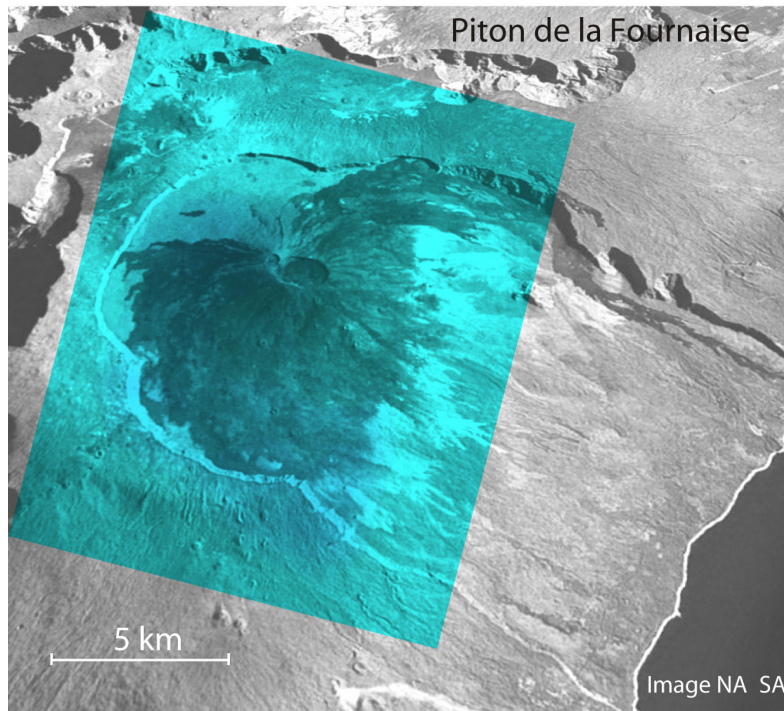This results in a **stretching of the waveforms**

### Synthetic velocity decrease



— reference CCF
— current CCF

time (s)

### Relative velocity variation measurement



$\Delta \tau / \tau = - \Delta \nu / \nu = 10\ \%$

moving window

$\tau$ (s)

$\Delta \tau$ (s)

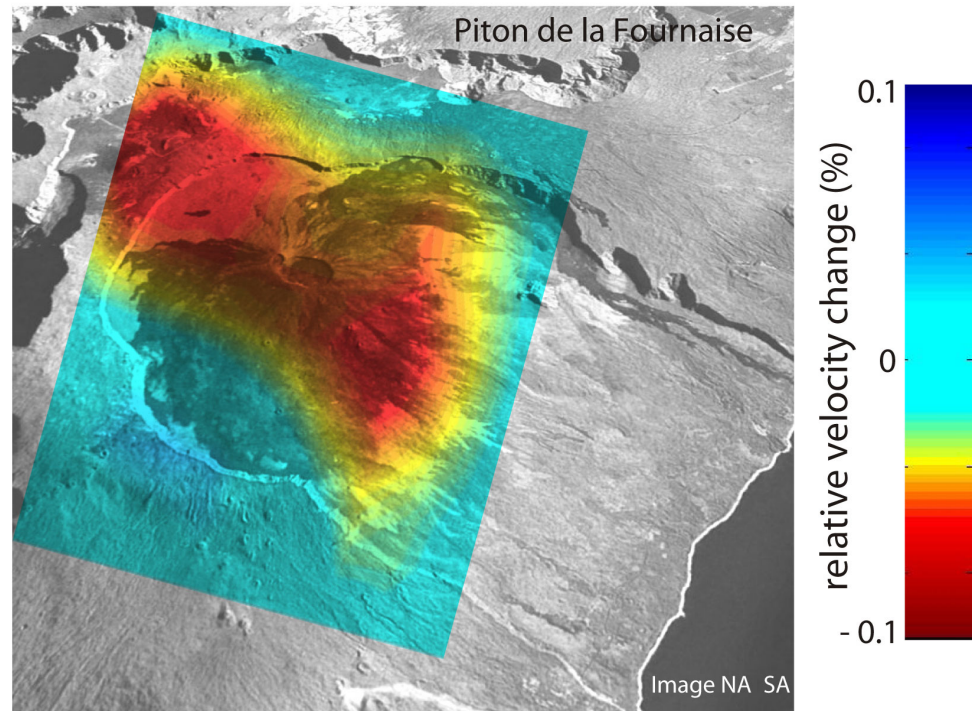# Monitoring velocity variations on a Volcano

# Monitoring velocity variations on a Volcano

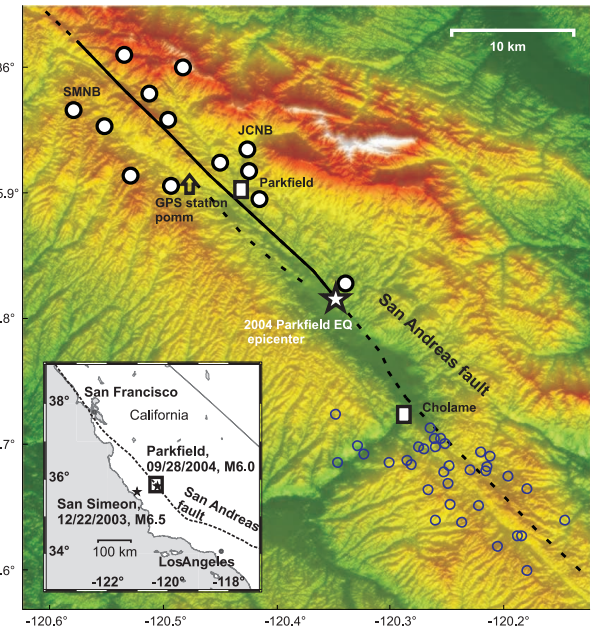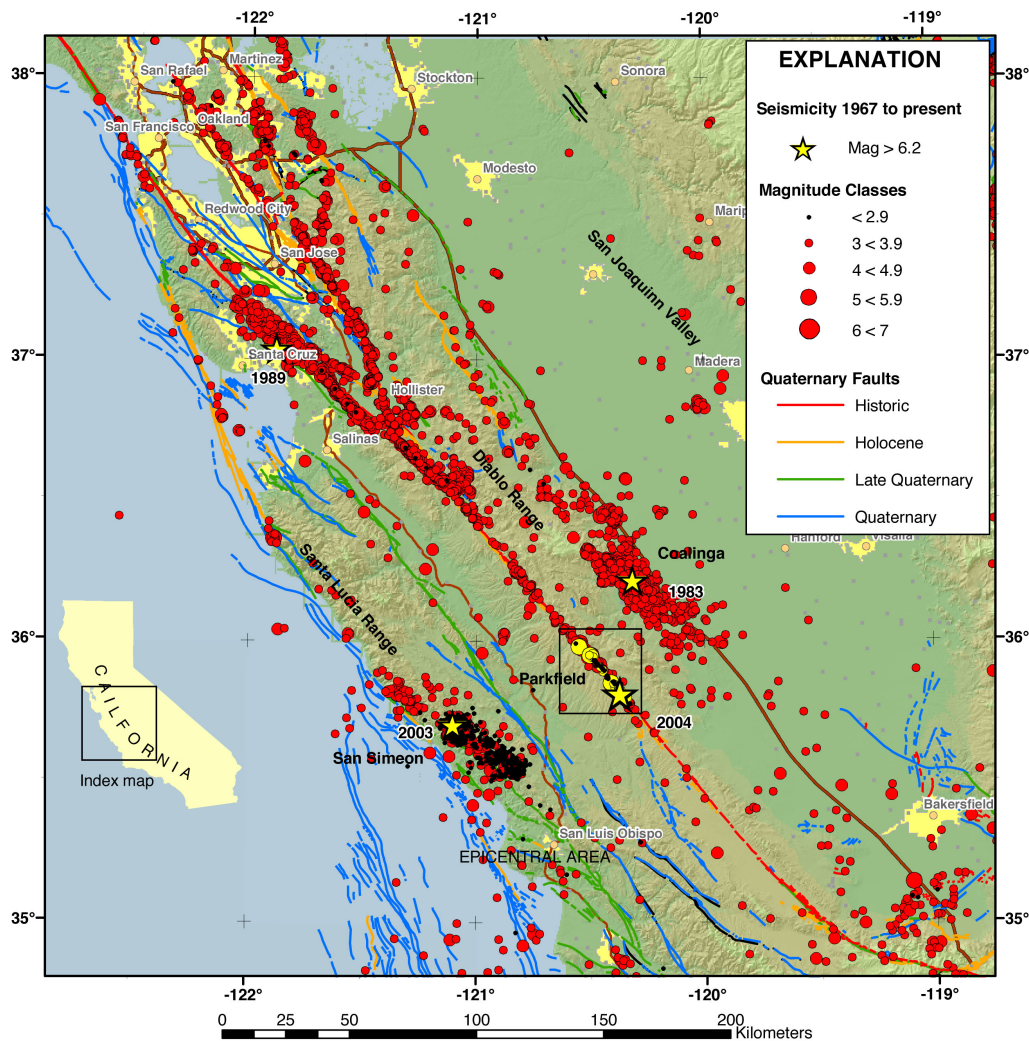**Short-term variations during 1999-2000: regionalization of the velocity perturbations**

9 days before eruption of June 2000

4 days before eruption of June 2000



Brenguier et al., 2008

# Monitoring velocity variations: San Andrea Fault



**PARKFIELD HIGH-RESOLUTION SEISMIC NETWORK**
**operated by Berkeley Seismological Laboratory**

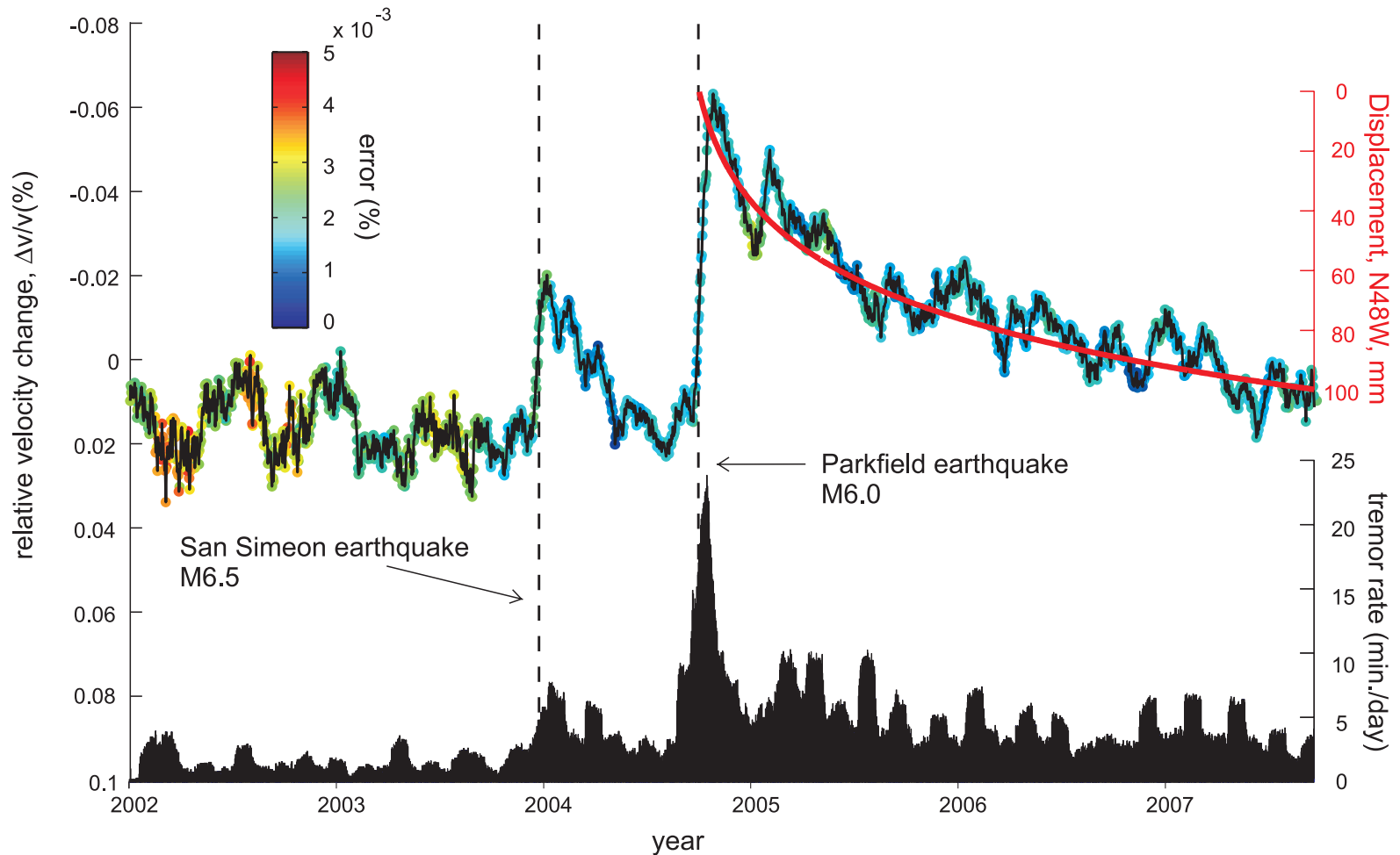Correlating and analyzing continuous seismic noise records during 2002-2007

0.1 – 0.9 Hz one day time window

Two M>6 earthquakes:

  M=6.6 San Simeon 2003 earthquake

  M=6.0 Parkfield 2004 earthquake

Brenguier et al. (2008)

Continuous excitation by oceanic gravity and infra gravity waves

Predominant peaks:

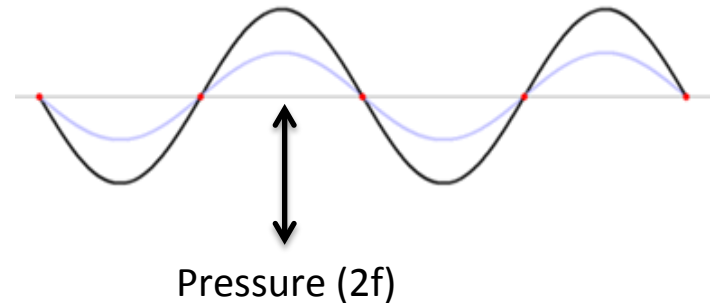- **Primary peak** : 10 – 20  s
- **Secondary peak**: 3 – 10s

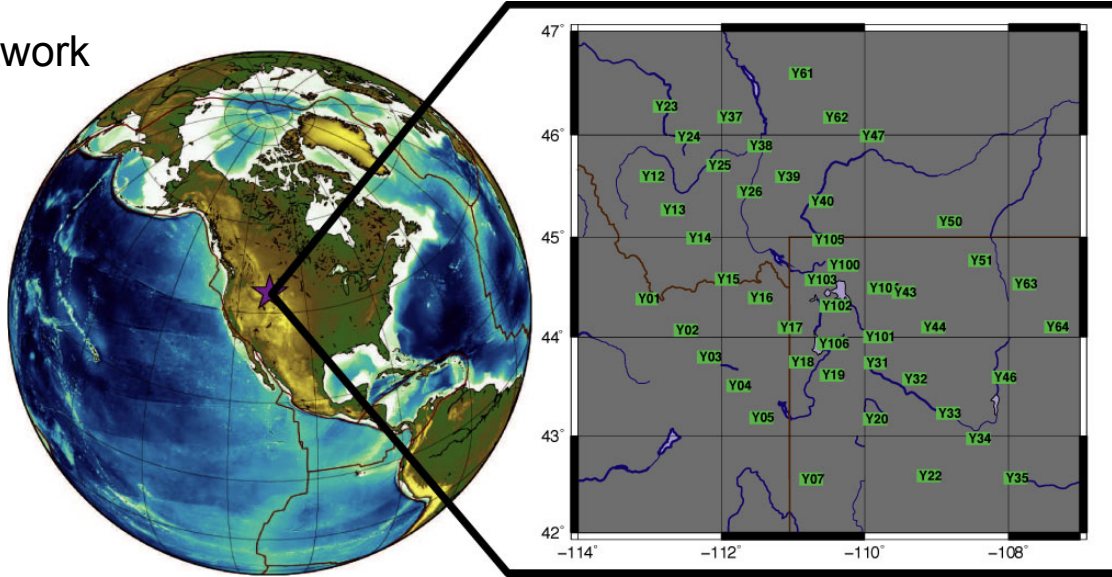Complex non linear interaction phenomena at coastlines and deep-sea oceans
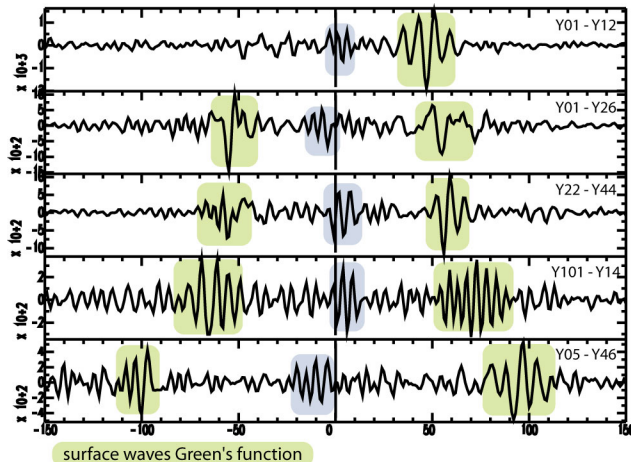


Ocean waves (f)        Ocean waves (f)

Pressure (2f)

# location of sources of the seismic noise can be investigated with processing continuous records of modern broadband seismic netorks and their correlations based on array based techniques
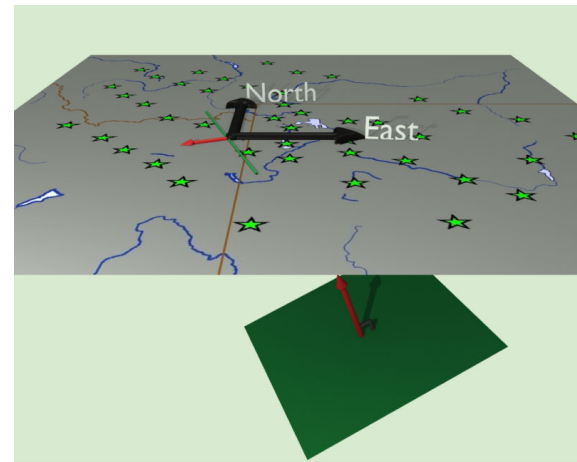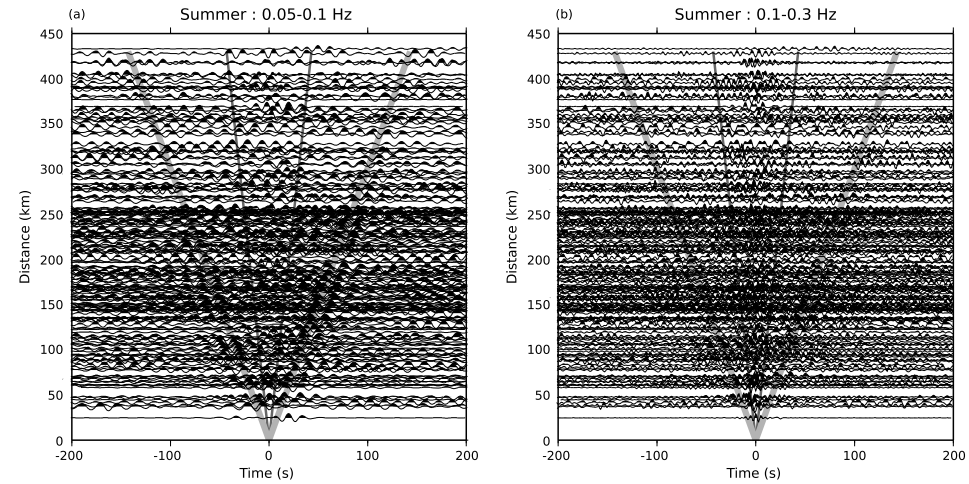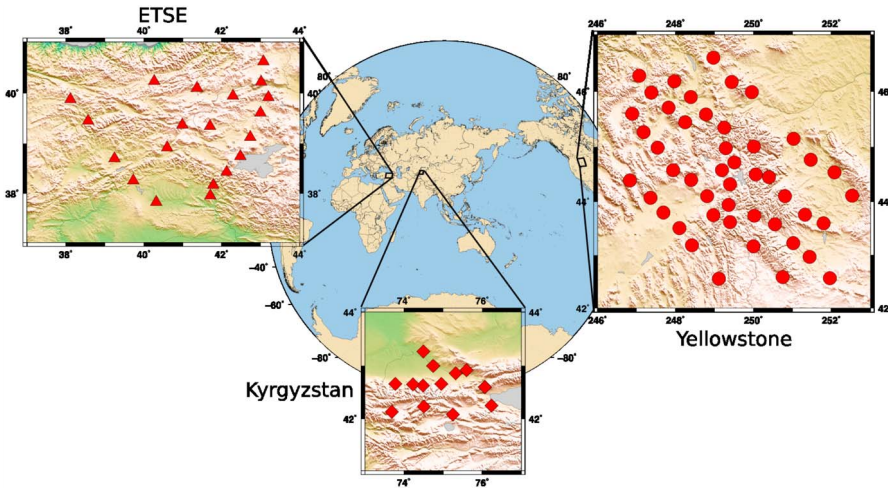
Yellowstone network



Noise correlations: clear arrivals at near-zero times - body waves from below



surface waves Green's function

Landes et al (2010)
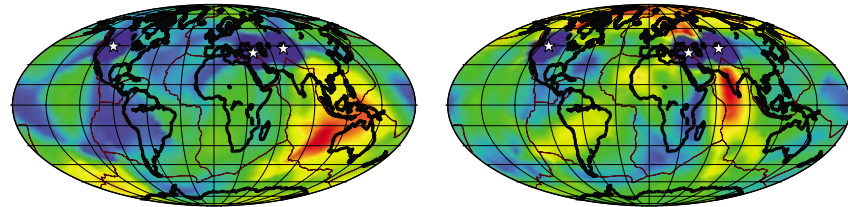
Noise sources are generated when there is interaction of ocean waves:

A. within a storm
B. by reflection at the coast
C. between storms

Source discretization:
Grid step=50km
Source= Vertical force at the ocean surface and random phase
Normal mode summation



Stutzman et al. (2012)
Gualtieri et al (2013)

Storm 1

A

B

Land

C

Class III

S2

Storm 2

Ocean waves are modelled every 6 hours
(code WAVEWATCH III, version 3.14,
6-hourly wind analysis from ECMWF).
Ocean wave interactions are computed
considering the 3 types of wave interactions
(Ardhuin et al. 2011)

Rayleigh waves    6 sec    Body waves

# Seismic noise correlation: Big Data



## Data ingestion / quality control

- N-dimensional *time series*
- *binary large objects (blob): > 100 TBs*
- *fine granularity: variable chunk sizes (GBs)*
- Partitioning, indexing, replication

## Data processing

- **Low level data access pattern**
- **Linear complexity**
- Streaming data workflow
- Provenance and metadata management

## Data analysis

- **Cross-correlation** and higher order statistics
- **Quadratic complexity** and CPU intensive
- Thread-blocks CUDA and CSP
- **Secondary data : ~ 6 * N² * $N_t$**
- Provenance and metada management

# Data-Intensive statistical analysis workflow



- Seismology PEs library and data streaming workflow (Dispel)
- Different execution models
- Data management layer: PFS
- Data management layer integration with value added analytics: iRODS platform + MonetDB
- Data provenance layer integration

# Data life cycles

**Persistent and resilient data**
- ✓ Public services for a wide community
- ✓ Data sets can range ~ 100 TB
- ✓ Hardware capacities and parallel capabilities

**Massive data processing pipelines**
- ✓ High bandwidth, optimal sequential IO and fast floating point operations
- ✓ Data volumes ~100s TB
- ✓ seismic noise correlation, image processing, high-rate GPS analysis
- ✓ Intermediate and derived data sets: ~100 TB
- ✓ *Lifecycle: weeks – months*

**Community analysis of very large data sets**
- ✓ Once massive data set arrives: partitioning and indexing, duplication
- ✓ Collaborative research data analysis and processing
- ✓ Scientific gateway, access policy, development environment
- ✓ Intermediate and derived data sets ~ 100 TB
- ✓ *Lifecycle: months -years*

# Data-intensive Infrastructure

## Intrinsic infrastructure mismatch

- Data volumes increase 100x in 10 years
- I/O bandwidth improves ~3x in 10 years
- Data analysis resources close to the data

## Need for efficient data crawling strategy

- **data locality**

  horizontal and vertical re-use

- **memory/IO bandwidth and latency**

  hierarchy of data storage (SDD,HDD)/memory, optimized aggregate sequential IO bandwidth

## Data Architecture:

- **Seismology database architecture:** archiving and distribution -> archiving synthetic models

- **Data processing architecture**: new data-intensive paradigms enabled by HPC, Hybrid architecture (GPU), PFS, HDFS, Hadoop-MapReduce; XLDB/MonetDB, CUDA-SQL, and MPI-DB toolkits

## A Data-scope environment and framework:

➢ **Analyze and model 100 TB+** of data in academic setting;

➢ At **least PB+ of storage with safe redundancy**;

➢ **High sequential IO throughput ~ aggregate disk speed**;

➢ **Streaming data analyses on par with data throughput**;

➢ **Distributed Infrastructures: HPC, Grid, Cloud**

## Infrastructure architecture:

- **A storage layer**: maximize capacity with enough disk bandwidth per server
- **A data-intensive processing layer:** maximize low level data access bandwidth and fabrics; fast sequential IO, large local disk storage, parallel file systems
- **A performance layer:** memory fabrics and bandwidth, CPGPUs, memory/disk hierarchy, interconnect bandwidth/latency
- **A development environment:** data and work flow engines with optimized data streaming, virtualization

# Data-intensive modelling: Earthquake Hazard assessment

*2001 Gujarati (M 7.7) Earthquake, India*



**Use parallel computing to simulate earthquakes and wave propagation (elastic/acoustic/ hydroacoustic)**

**Learn about structure of the Earth based upon seismic waves (tomography)**

**Produce seismic hazard maps (local/regional scale) e.g. Los Angeles, Tokyo, Mexico City**
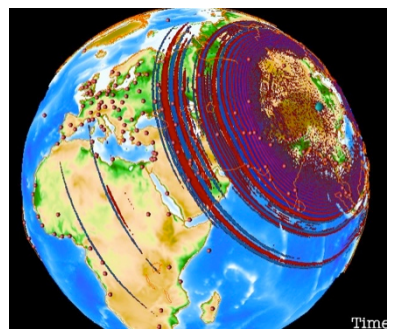
20,000 people killed
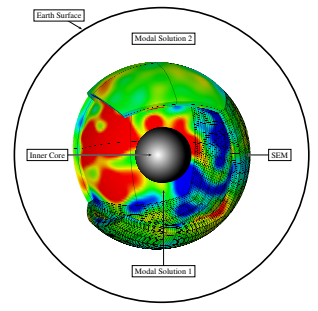167,000 injured
≈ 339,000 buildings destroyed
783,000 buildings damaged

# Data-intensive HPC simulation

**Seismic wave propagation**


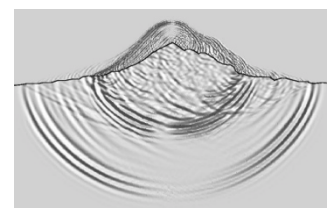
Komatisch *et al.* (2009)



Capdeville *et al.* (2003)

**Global scale:**

- **Waveform prediction for large earthquakes**
- **Understanding complex wave propagation at global scale**

**Aero-acoustic wave simulation in a volcano**





Käser *et al.* (2009)

**Regional scale:**
- **Waveform prediction in complex media**
- **Seismic/acoustic/Hydroacoustic coupling**

**Strong motion simulation:  Grenoble Valley**



Chaljub *et al.* (2009); Delavaud *et al.* (2009), Käser *et al.* (2009)

**Strong motion prediction:**
- **Physically-based hazard assessment**
- **Earthquake source dynamics**
- **Stochastic earthquake scenarios**
- **Stochastic wave simulation**

# Specfem3D: a community code





**Goal: model acoustic / elastic / viscoelastic / poroelastic / seismic wave propagation in the Earth (earthquakes, oil industry), in ocean acoustics, in non destructive testing, in medical acoustic tomography…**

The SPECFEM3D source code is open (GNU GPL v2)

Mostly developed by Dimitri Komatitsch and Jeroen Tromp since 1996.

Improved with INRIA (Pau, France), CNRS (Marseille, France), the Barcelona Supercomputing Center (Spain) and University of Basel (Switzerland).

# Variational Formulation: Solid case

Differential or *strong* form:

$$\rho\, \partial_t^2 \mathbf{u} = \nabla \cdot \sigma + \mathbf{f}$$

Variational or *weak* form in the time domain:

$$\int \rho\, \mathbf{w} \cdot \partial_t^2 \mathbf{u}\, \mathrm{d}^3\mathbf{r} = -\int \nabla \mathbf{w} : \sigma\, \mathrm{d}^3\mathbf{r}$$

$$+ \int \nabla \mathbf{w} : M(\mathbf{r}_s)\, S(t)\, \mathrm{d}^3\mathbf{r} - \int_S \mathbf{w} \cdot \sigma \cdot \hat{\mathbf{n}}\, \mathrm{d}^2\mathbf{r}$$

+ attenuation (memory variables) and ocean load

# Spectral Element Method

- Accuracy of a spectral method, flexibility of a finite-element method

- Extended by Vilotte, Komatitsch, Capdeville, Chaljub, Tromp…

- "spectral" finite-elements with high-degree polynomial interpolation

- Gauss-Lobatto-Legendre quadrature

- Explicit high-order time integration

- Very efficient on parallel computers, no linear system to invert (diagonal mass matrix)

- Can be extended through a high-order Discrete Galerkin approximation

# Porting Specfem3D on GPU

- At each iteration of the serial time loop, three main types of operations are performed:

  - **update (with no dependency) of some global arrays composed of the unique points of the mesh**

  - **purely local calculations of the product of predefined derivative matrices with a local copy of the displacement vector along cut planes in the three directions (i, j and k) of a 3D spectral element**

  - **update (with no dependency) of other global arrays composed of the unique points of the mesh**
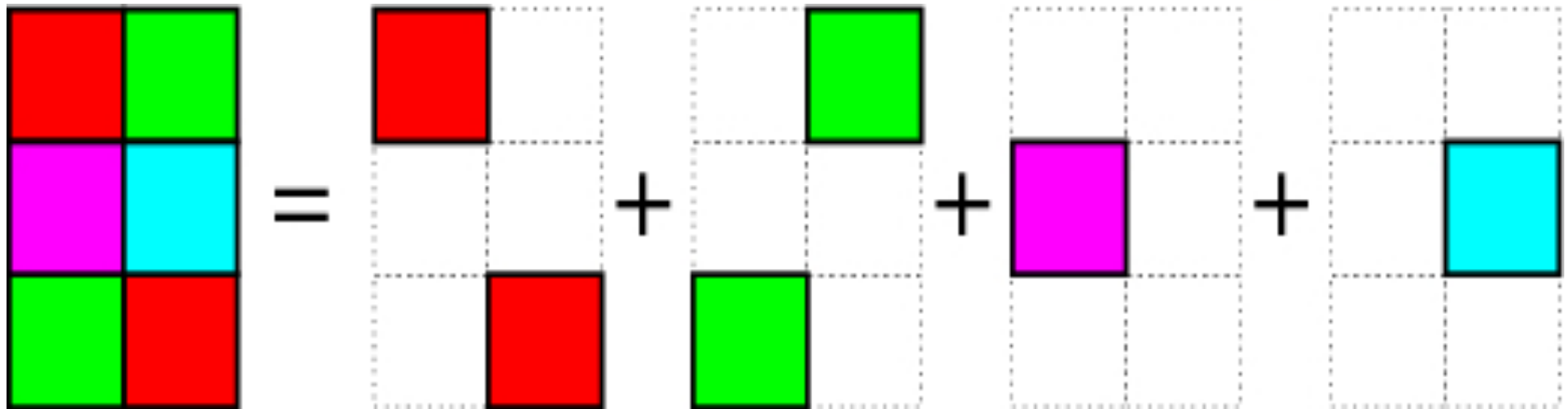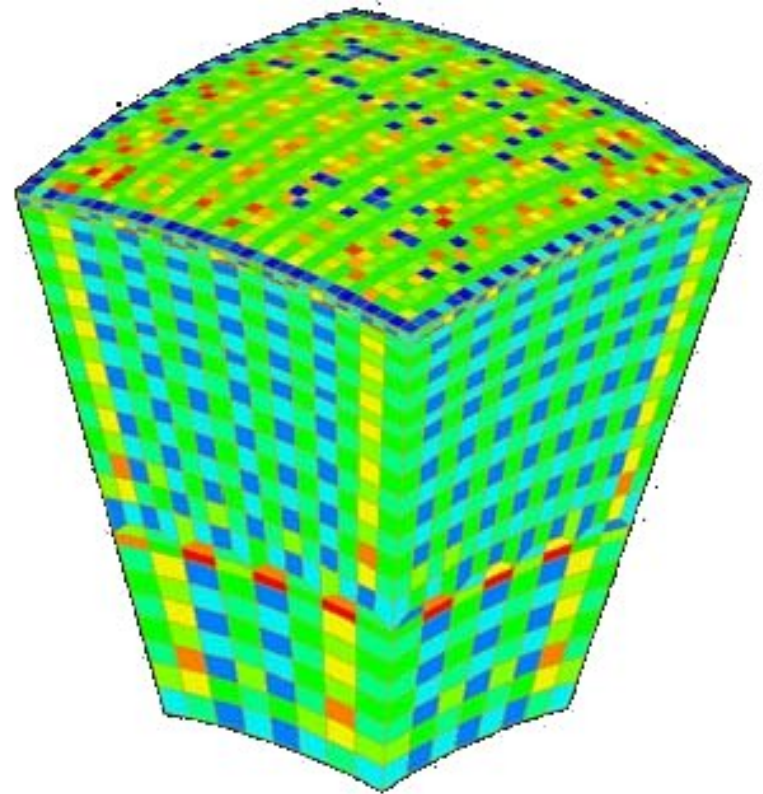
# Minimize CPU/GPU data transfers

- CPU ↔ GPU memory bandwidth much lower than GPU memory bandwidth

- Use page-locked host memory (cudaMallocHost()) for maximum CPU ↔ GPU bandwidth

- Minimize CPU ↔ GPU data transfers by moving more code from CPU to GPU, even if that means running kernels with low parallelism computations

- Intermediate data structures can be allocated, operated on, and deallocated without ever copying them to CPU memory

- Group data transfers: one large transfer much better than many small ones

- Fit all the arrays on the GPU card to avoid costly CPU ↔ GPU data transfers

- But of course the MPI buffers must remain on the CPU, therefore we cannot avoid a small number of transfers (of 2D cut planes)
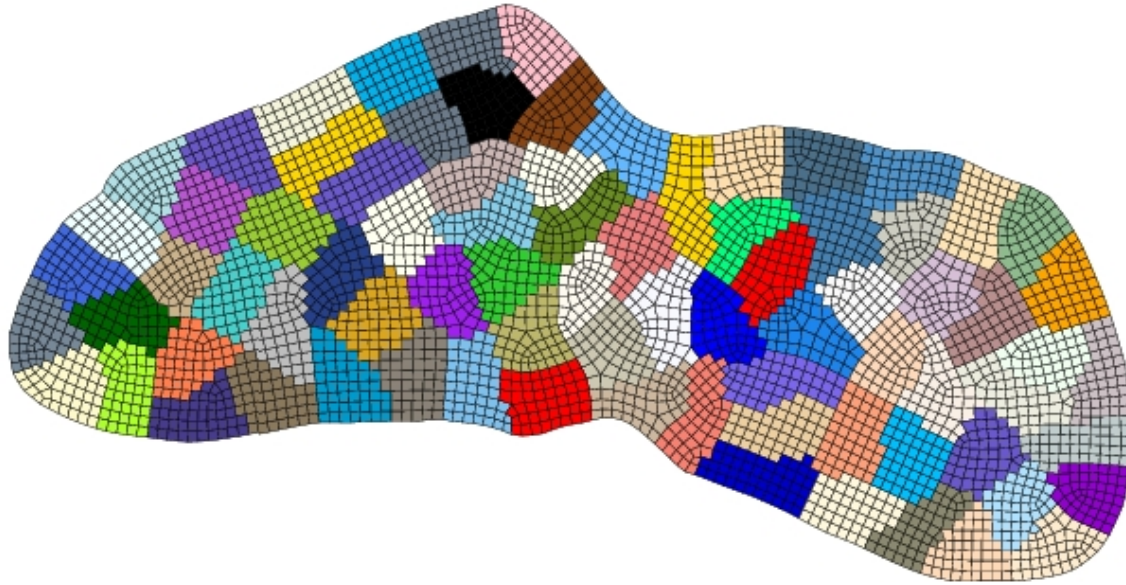
# Mesh Coloring

Ensure that contributions from two local nodes never update the same global value from different warps

Use of mesh coloring: suppress dependencies between mesh points inside a given kernel
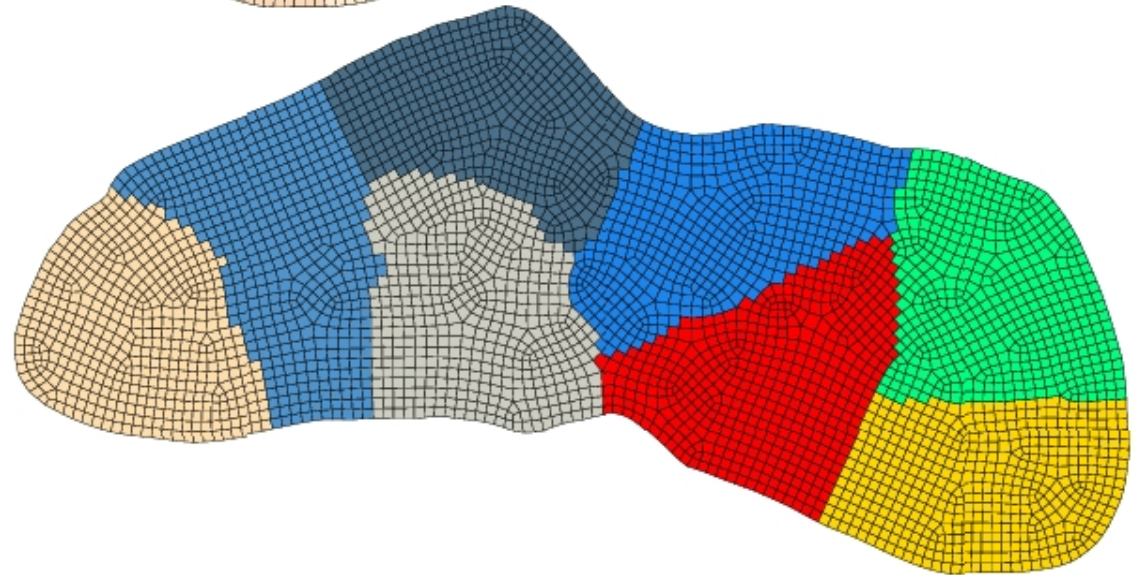
Use of "atomic" leads to slower code

# Non Blocking MPI to overlap



80 domains : the inner part is too small to overlap MPI communications or CUDA data transfers with calculations.

**Danielson and Namburu (1998)**

8 domains: granularity is good and we can overlap.

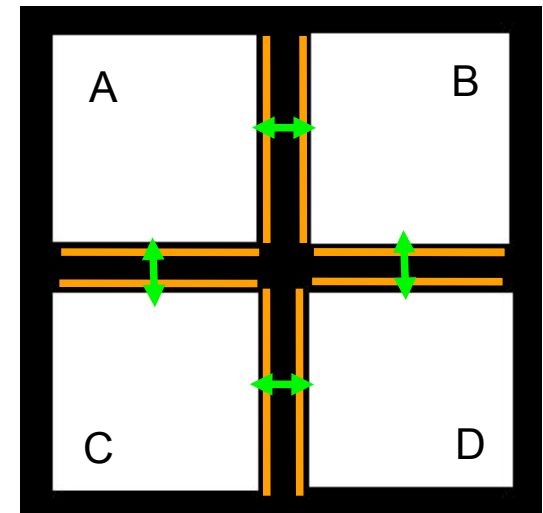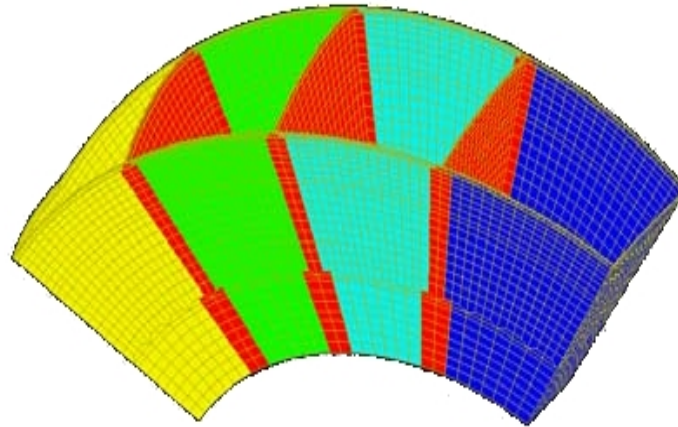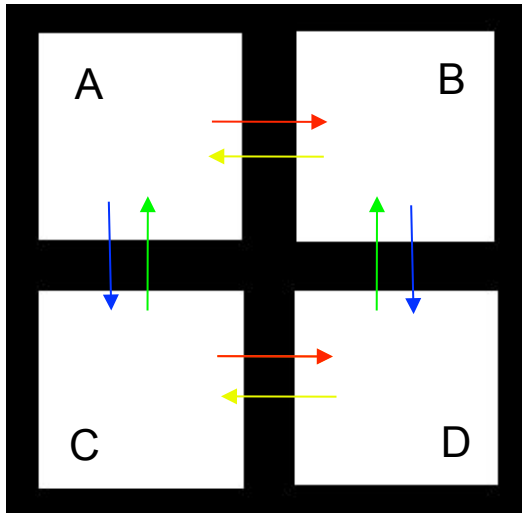**D. Komatitsch in collaboration with Roland Martin and Nicolas Le Goff (INRIA, Pau, France)**

# Adding MPI to GPU

- Old communication scheme (blocking MPI)
- Update done in the whole arrays (all elements computed before starting MPI calls)

New communication scheme (non blocking MPI)

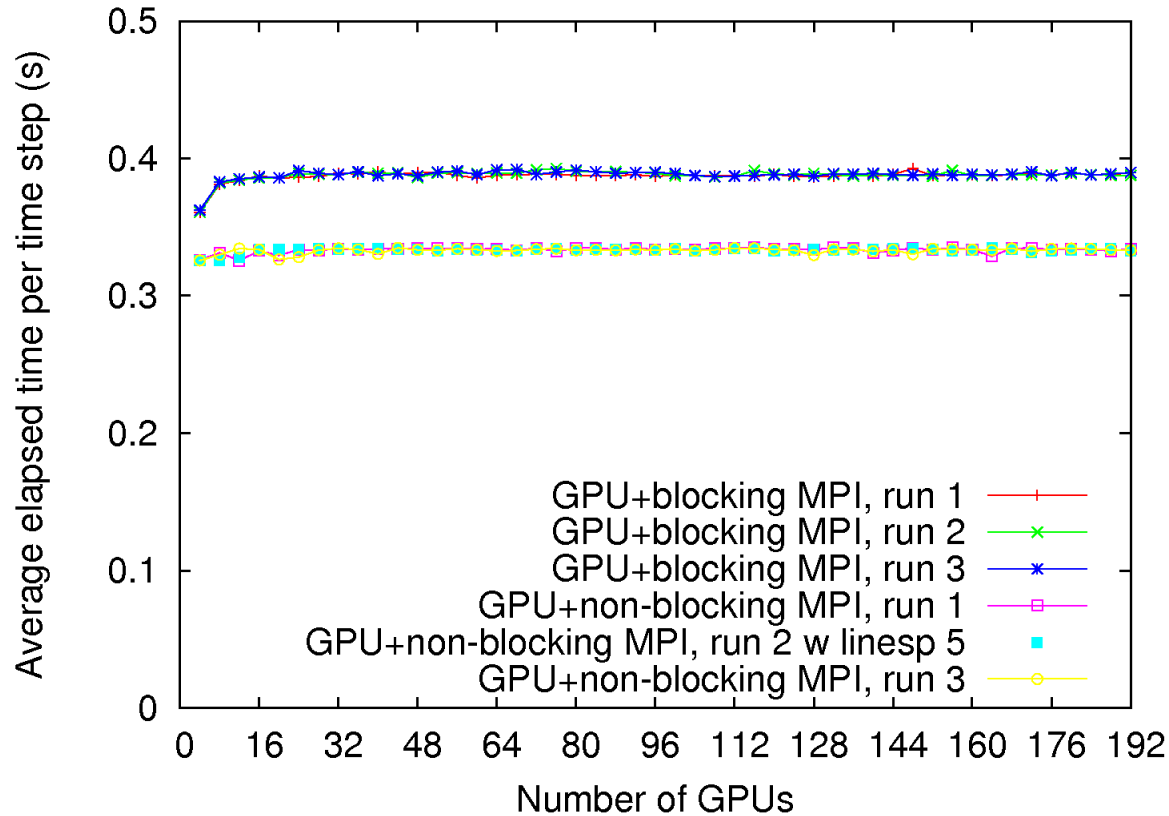Update done in buffers (for outer mesh elements first)



MPI communications cost on GPU version ~ 5%,

> We need to use non-blocking MPI communications.

> MPI communications are very well overlapped by computations on the GPU.

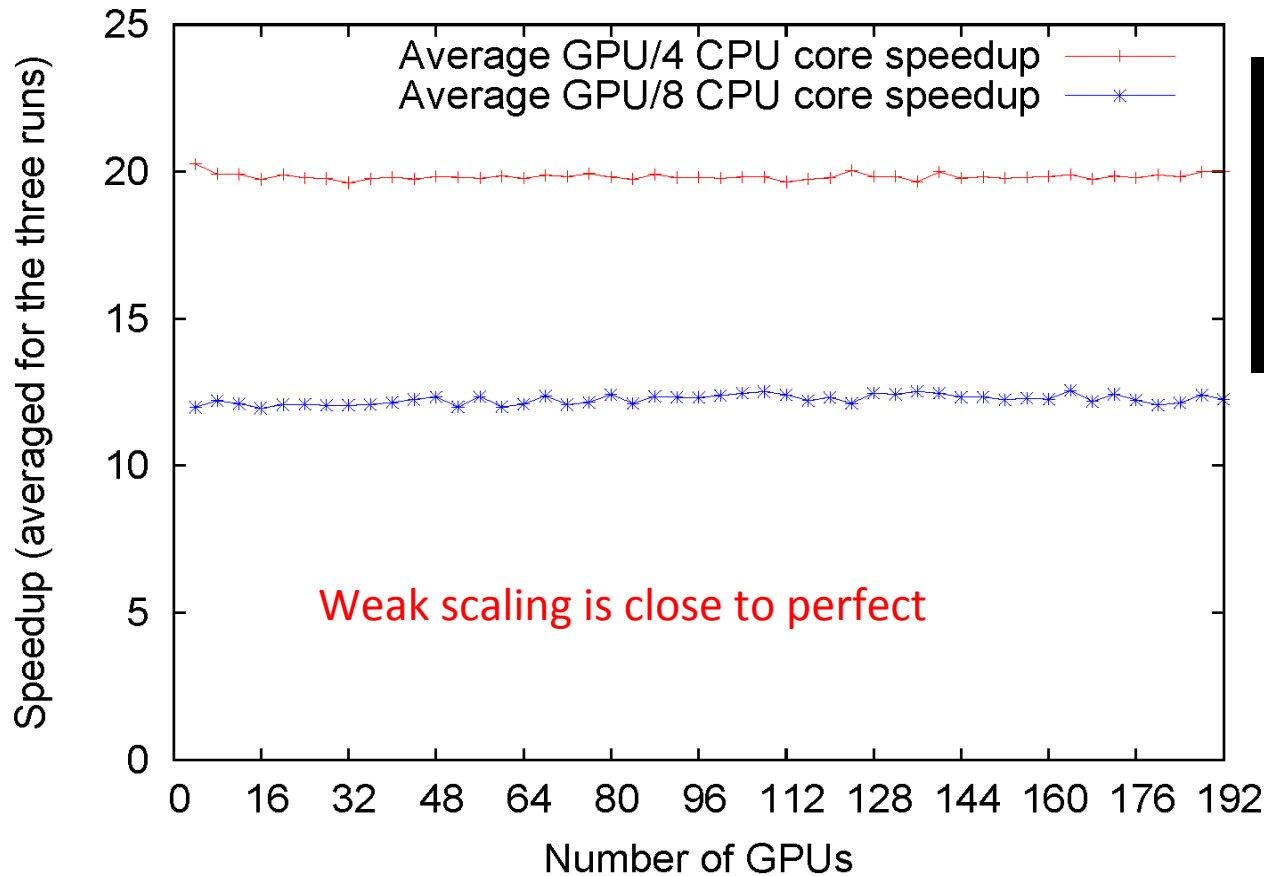# MultiGPU weak scaling (up to 192 GPUs)



- Constant problem size of 3.6 GB per GPU
- Weak scaling excellent up to 17 billion unknowns
- Blocking MPI results in 20% slowdown

It is difficult to define speedup: versus what?

On the CEA/CCRT/GENCI GPU/Nehalem cluster, **about 12x versus all the CPU cores, 20x for one GPU versus one CPU core.**
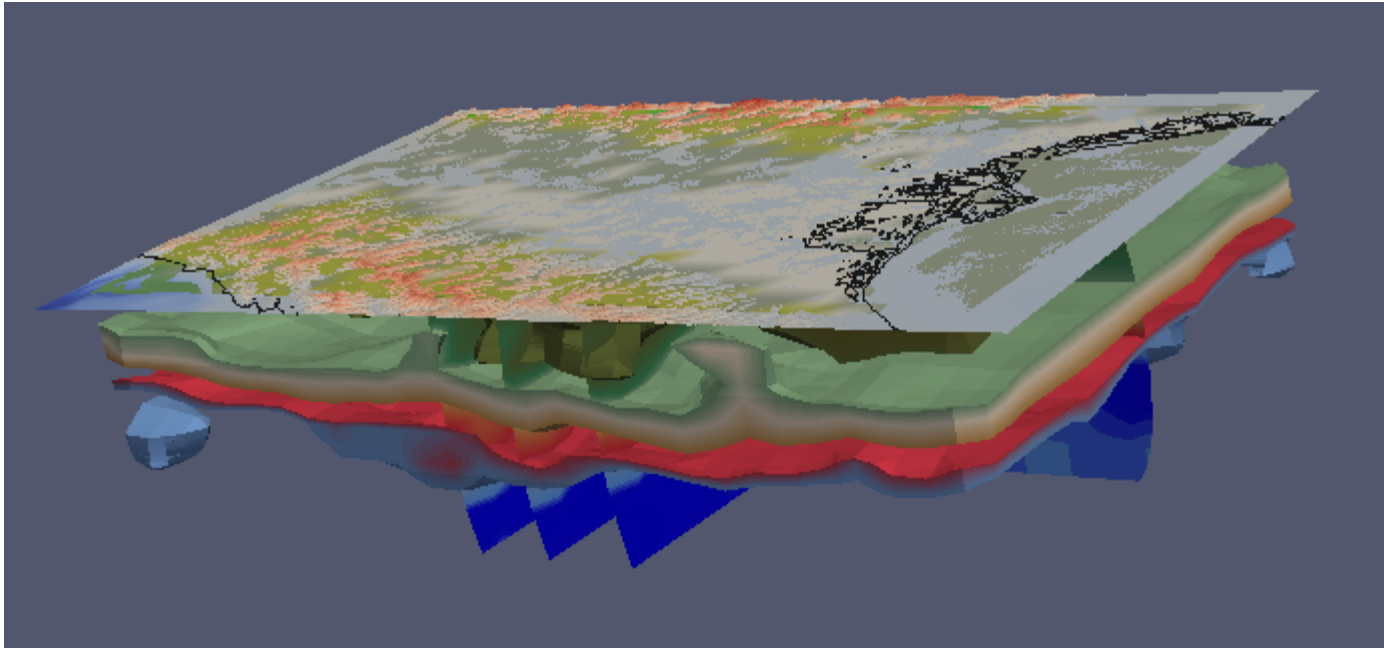
# Multi-GPU weak scaling (up to 192 GPUs)



High-frequency ocean acoustics, inverse problems in seismology, acoustic tomography, reverse-time migration in seismics: high resolution needed, and/or large iterative problems to solve ⇒ Large calculations to perform.
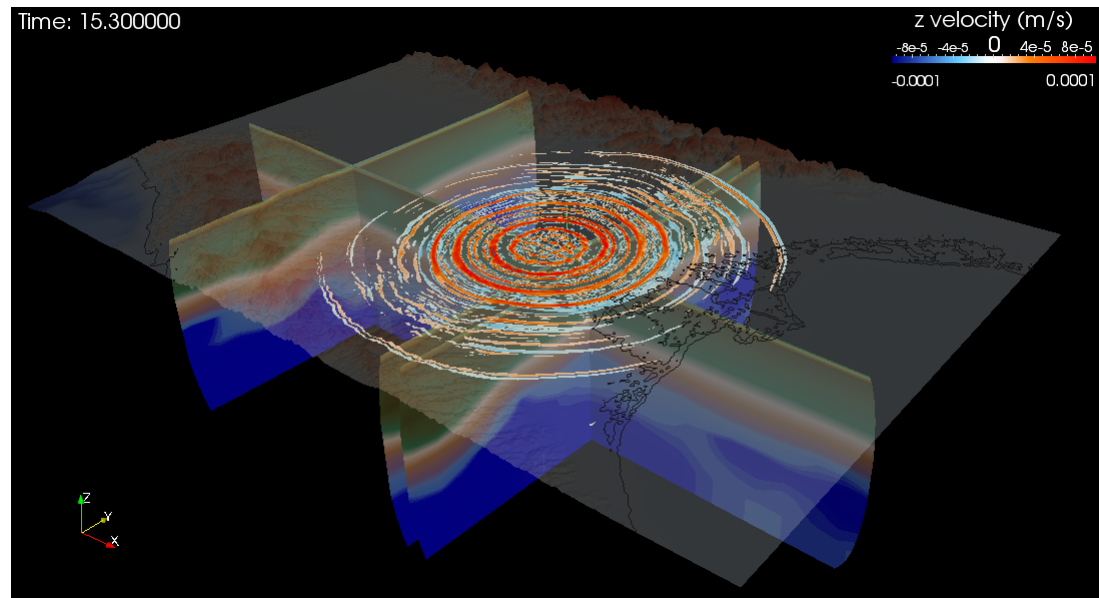
⇒ GPU computing: code needs to be rewritten, but large speedup can be obtained (around 20x-30x for Specfem3D, but it is difficult to define speedup).

# Northern Italy event of May 20, 2012



**Collaboration D. Komatitsch
with INGV (Emanuele
Casarotti et al., Roma and
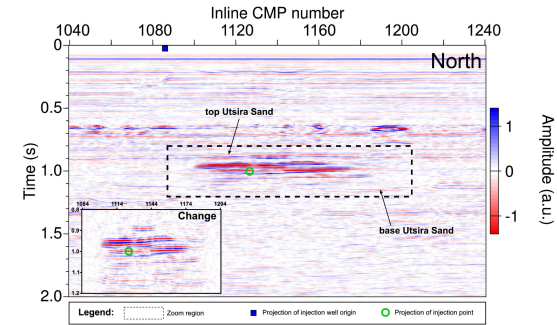Irene Molinari et al., Bologna)
+ CASPUR + CINECA.**

**Run on CASPUR machines.**

# Data-intensive modelling: adjoint-based inversion

**Marine exploration geophysics:**
**High Resolution Imaging (inversion/migration)**

**Exploration geophysics: Time lapse HR tomography imaging**

**Full wave form Tomography: Global scale**
**Unrevealing the Earth's structure**

**HR regional Tomography/migration:**
**Unrevealing subduction structure**

# Data-intensive modelling: Full waveform inversion

## Adjoint Tomography of Europe

### "Big Data"

| earthquakes | stations | iterations | simulations | CPU hours | measurements |
|---|---|---|---|---|---|
| 190 | 745 | 30 | 17,100 | 2.3 million | 123,205 |

one day of seismic record

noise

coda

ballistic waves used in traditional tomography

Krishnan et al (2012)

Tromp et al (2012)

Depth 75 km

Hejun Zhu

# Adjoint-based methods

$$\chi_1(\mathbf{m}) = \frac{1}{2}\sum_{r=1}^{N_r}\int_0^T w_r(t)||\mathbf{s}(\mathbf{x}_r,t;\mathbf{m}) - \mathbf{d}(\mathbf{x}_r,t)||^2\,\mathrm{d}t,$$

$$\delta\chi_1 = \int_V [K_\rho(\mathbf{x})\,\delta\ln\rho(\mathbf{x}) + K_\mu(\mathbf{x})\,\delta\ln\mu(\mathbf{x}) + K_\kappa(\mathbf{x})\,\delta\ln\kappa(\mathbf{x})]\,\mathrm{d}^3\mathbf{x},$$

$$K_\kappa(\mathbf{x}) = -\int_0^T \kappa(\mathbf{x})\,[\nabla\cdot\mathbf{s}^\dagger(\mathbf{x},T-t)]\,[\nabla\cdot\mathbf{s}(\mathbf{x},t)]\,\mathrm{d}t,$$

<u>Theory</u>: A. Tarantola, Talagrand and Courtier, Virieux, Singh, Tromp.

Close to time reversal (Mathias Fink et al.) but not identical,
        thus interesting developments to do.

# CPU-intensive modelling: waveform inversion



## High performance parallel codes

- Specfem3D, Seisol ...

## Waveform inversion

- Non-linear inversion
- Adjoint-based inversion methods: **->** one forward and one adjoint simulations per Newton iteration for each time step and earthquake

## Orchestrated workflow
- Data Intensive analysis and High Performance computing
- Across Public HPC and Private data and computing infrastructures

## Big Data
- Earthquake event waveforms: synthetics and observed
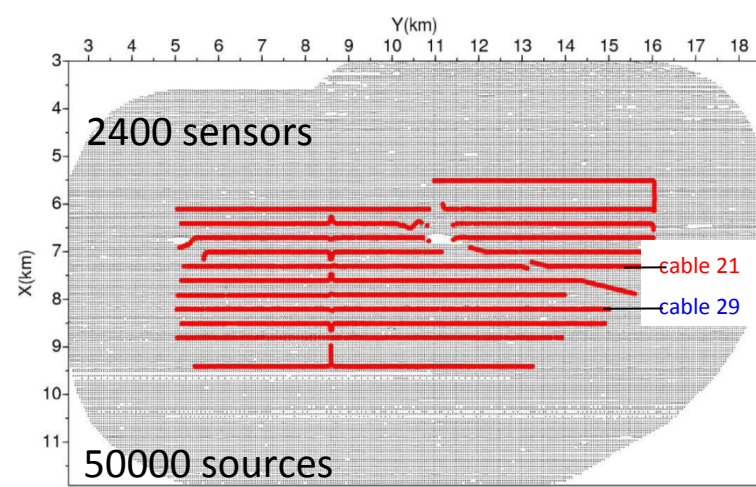- State of the systems: x,y,z,t -> v, σ

## Mesh generation

- Quality control and parallel mesh generation

3D real application:
Valhall case

Virieux and collaborators
(Seiscope)

2400 sensors

cable 21
cable 29

50000 sources

- Exploited since 1982, Life of Field Seismic (LoFS) network since 2003

- BP starting models by anisotropic reflection traveltime tomography

- Strong imprint of anisotropy in the seismic Valhall dataset

- 3D isotropic acoustic FWI by Sirgue & al. (2010) using 13 km maximum offset

*(Sirgue & al., 2010)*



cable 21

cable 29

velocity (m/s)

Z=1050m

velocity (m/s)

# 3D acoustic FWI

Brossier et al (2013)

| For+Inv | Few cores | Many cores |
|---------|-----------|------------|
| Time+Freq | 20830 s | 326 s |
| Freq+Freq | 6209 s | 1445 s |

3D monoparametric reconstruction
(Pratt's strategy)

Etienne et al (2012), Hu et al (2012)

# 3D acoustic FWI



Result at 4 Hz - Horizontal cross sections

Superficial channels

Gas reservoir

Imprint of the acquisition

Vp FWI

# 3D acoustic FWI



Result at 7 Hz - Horizontal cross sections

Superficial channels

Gas reservoir

Imprint of the acquisition

# Data-intensive HPC workflow



- Orchestrated workflows and execution models
- Stream based data analysis and enabled CSP wave simulation codes (Specfem3D and Seisol)
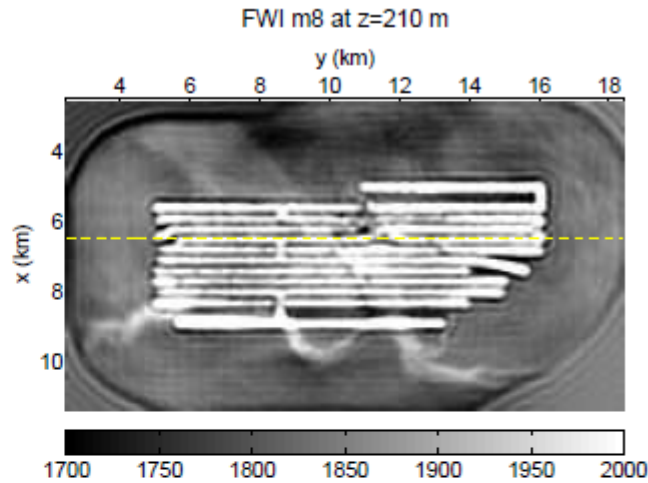- Job submission across Grid & HPC DCIs: AAA (X.509 proxies), JSAGA/DCI-Bridge
- Data streaming and files transfer orchestration across DCIs:
- GridFTP enabled data transfer PEs, iRODS

# PRACE Infrastructure: French TGCC





~2 Petaflops for the European infrastructure



**The TGCC (Très Grand Centre de Calcul / "Very Big Computing Center")** hosts the PRACE "CURIE" European machine

GENCI in France, CINECA / CASPUR in Italy.

## Data-intensive computing challenges

**Large scale 3D simulation**:

- multi-scale and multi-physics
- stochastic direct uncertainty evaluation

**Inversion and Data assimilation**:

- adjoint-based methods: non linear iterations with large number of forward and adjoint simulations
- stochastic methods: inverse uncertainty quantification

**Orchestrated workflows**:

- data analysis and modeling applications
- end-to-end applications

**Hilbert SFC of level 2 and 64 sub-cubes**



**Domain decomposition by METIS (left) and SFC (right)**

### Scalability

Communication fabrics
Asynchronous time integration, vertical reuse
Explicit locality model (vertical/horizontal)
Parallel large system solver
Dynamic load balancing

### Data-intensive HPC

Memory hierarchy and bandwith
Fast sequential IO
Hierarchy of storage HDD/SDD
Advanced data-structure and parallel filesystems

### Multicore architectures

Mixed-hybrid parallel implementation
High-level task concurrency: asynchronous task parallelism; overlapping computation and communication
Self-scheduling at task level
Fault tolerance system

### End-to-end analysis

Parallel unstructured mesh generation
Domain decomposition
Post-processing data-intensive data analysis
Data management

# A service-oriented architecture

**Separation of concerns**

**Workbenches for seismologists**

Iterative
data-intensive
development of
research
methods

Abstract
level

Registry

Gateway interface
one integrated
model

Enactment
level

Mapping
optimisation
and
distributed enactment

Accommodating
**Many groups of researchers**
**Many tool sets**
**Many research strategies**
**Many working practices**

**Canonical representation**

Composing or hiding
**Many autonomous resources & services**
**Multiple enactment mechanisms**
**Multiple platform implementations**
**Multiple e-Infrastructures**

**System heterogeneity and complexity**

**Resilience toward "standards" evolution**

# Architecture

**Architectural changes**

- Tipping balance to data : data crawling architecture strategy;
- Support both Big Data DC architectures: data-intensive analysis – loosely coupled, data streaming on par with data throughput - and CPU-intensive architecture – tightly coupled;
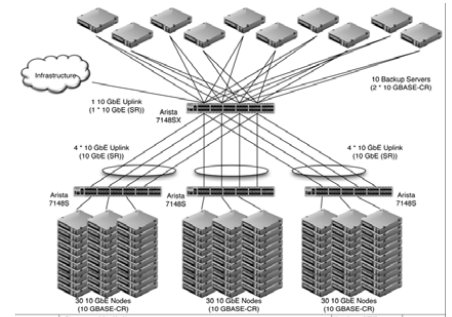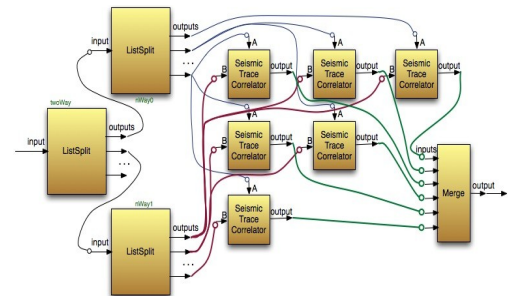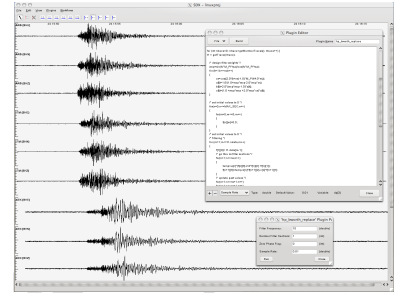- Compute in storage architecture and technology with added analytics;
- Augmented hierarchical object-based storage management, and heavy concurrent data access beyond POSIX;

**What operational changes**

- Supporting extended Data life-cycle within HPC infrastructures: data storage hierarchies and scientific gateways;
- Analytics platform must integrate Data-intensive HPC infrastructures and Data-intensive HTC infrastructures;
- Supporting orchestrated workflow – and data flow - across BD and EC DCIs and execution models: access policy, AAA mechanism, monitoring tools ….

# Software

**Data management/exploration**

- PFSs, iRODS, Scientific data bases (MonentDB)
- Data archives: Data and Metadata structure (<- acquisition/transmission & data exchange format)

**Software library and tools**

- Analysis domain specific libraries: ObsPy, Python, NumPy, SciPy, SeisHub, C/C++, Matlab
- 3D wave simulation codes (Specfem3D and Seisol) continuous optimization. Good strong and weak scaling up to ~30-40 K cores.

**Data management system needs**

- Beyond Posix : n-dimensional objects, Blobs with dynamical adjustable chunk size, storage; concurrent access, versioning-based concurrent access
- Explore self-describing formats: HDFS, NetCDF, ADIOS

**Software missing**

- Fault tolerance: workflow & HPC codes (FTI experiments with Specfem3D, Bautista-Gomez et al., 2011)

**Big Data**

**Data Archives and Data infrastructure**

Global observation systems: Integrated distributed
data archives
Long term observatories: raw data preservation,
data curation, data annotation
Data and Metadata standards
Data management and data exchange standards

**Data-intensive research**

- Increasingly large data sets (> 100-500 TBs each)
- Data-intensive: HPC modelling (inversion/assimilation); statistical analysis
- Different data life cycle:
    - ➢ Long-term (years) with shared services;
    - ➢ Mid-term (1-2 years), for research group analysis/modelling;
    - ➢ Short-term (few months) for massive processing (on demand ?) pipelines.
- Hierarchy of distributed storage -> vertical reuse optimization
- Orchestrated workflow across HPC infrastructures and Grid-like private/public infrastructures
- Secondary products publish in the Data archives with provenance and metadata
- Continuous data curation process