# A BRIEF HISTORY OF THE SSDBM CONFERENCE SERIES

## 30TH ANNIVERSARY

## Arie Shoshani

### Lawrence Berkeley National Laboratory

SSDBM conference
July 9-11, 2018

# Outline

- How did this conference series start
- Research topics evolution over time
- Future challenges
- Light-hearted anecdotes
- Next conference – Santa Cruz, California

# 30 SSDBM conferences over 37 years

## PREVIOUS CONFERENCES

2018, Bozen-Bolzano, Italy
2017, Chicago, Illinois
2016, Budapest, Hungary
2015, San Diego, California
2014, Denmark
2013, Baltimore
2012, Crete, Greece
2011, Portland, Oregon
2010, Heidelberg, Germany
2009, New Orleans
2008, Hong Kong
2007, Banff, Canada
2006, Vienna, Austria
2005, Santa Barbara, California
2004, Santorini, Greece
2003, Cambridge, Massachusetts
2002, Edinburgh, Scotland
2001, Fairfax, Virginia
2000, Berlin, Germany
1999, Cleveland, Ohio
1998, Capri, Italy
1997, Olympia, Washington
1996, Stockholm, Sweden
1994, Charlottesville, Virginia
1992, Ascona, Switzerland
1990, Charlotte, North Carolina
1988, Rome, Italy
1986, Luxembourg
1983, Los Altos, California
1981, Menlo Park, California

## OBSERVATIONS

- Great locations

- Great social experience

- Small crowd, no parallel sessions

- All volunteer work

- Based on popular interest


- I attended all, but one

- I had papers in most


- Next: Santa Cruz, California

A. Shoshani

# Department of Energy Labs



Office of Science Labs
Other Offices Labs

A. Shoshani

# DOE's Leadership Class Facilities

**Oak Ridge Leadership Computing Facility**
Titan
Cray XK7
20 petaflops
hybrid-architecture
18,688 AMD 16-core Opteron 6274 CPUs (a total of 299,008 processing cores)
18,688 NVIDIA Kepler GPUs
710 terabytes of memory
10 petabyte disk

**Argonne Leadership Computing Facility**
Mira
IBM Blue Gene/Q
10 petaflops
786,432 processors
768 terabytes of memory
7.6 petabytes disk

**NERSC** The National Energy Research Scientific Computing Center (NERSC) - LBNL
Hopper
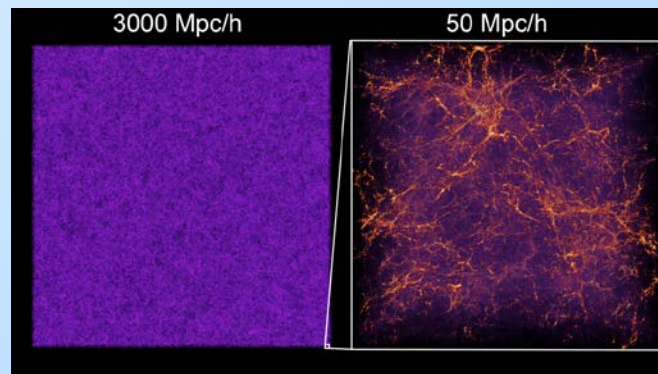Cray XE6
1.28 Petaflops/sec,
153,216 compute cores,
212 Terabytes of memory, and
2 Petabytes of disk.

**ESnet**
Energy Sciences Network (ESnet)
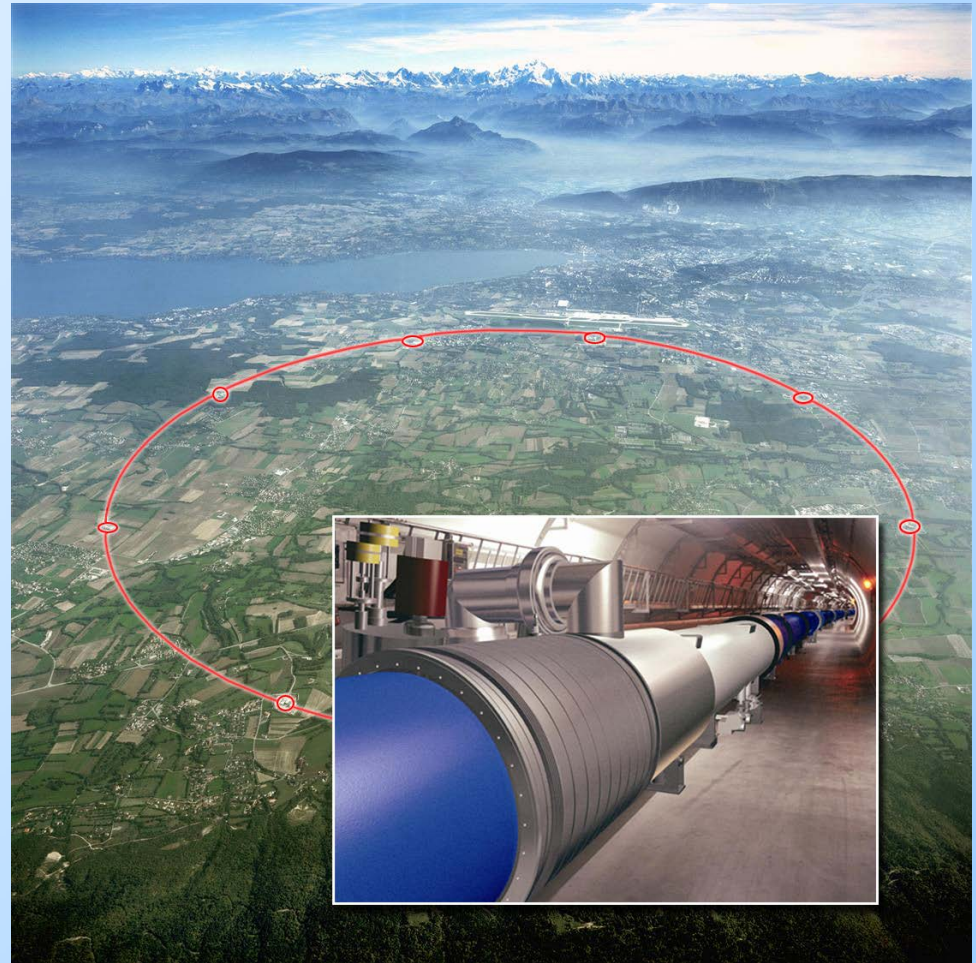Upgraded recently to 100 Gb/s on main connections



A. Shoshani

# Example of Large Data Volume in Science

## Large Hadron Collider: to find the God particle

- sensors capable of 140PB/s
- reduce 99.99% of data by hardware triggers
- Keep 15 PB per year
- 27 km tunnel
- ~10,000 superconducting magnets
- Operating temperature 1.9 Kelvin
- Construction cost: US$9Billion
- Power consumption: ~120 MW

A. Shoshani

# Data models and SSDBM

- Pre-1970
  - Hierarchical model
    - [Integrated Data Store](#) (IDS), by GE
    - Model based on efficient physical organization
      - E.g. projects $\longrightarrow$ employees, employee $\longrightarrow$ children
      - Specialized query interfaces (procedural: follow pointers)
    - Later: XML databases
    - Problem: data model does not capture more complex associations: projects $\longleftrightarrow$ employees
- Post-1970
  - Relational model
    - Separation of logical data model from physical data model (physical data independence)
    - Logical-level query language (SQL)
    - Mapping required query optimization, indexing, physical data layout,
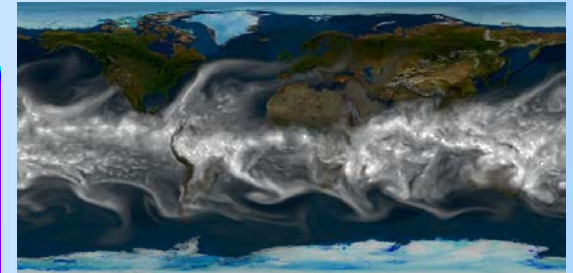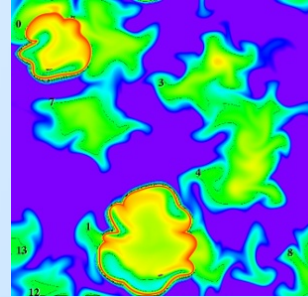    - Multiple implementation based on a standard query language

A. Shoshani

# Why Scientists Don't Use Data Management Systems?

## (when I Joined LBNL in 1976)

A. Shoshani

# What does "Scientific Data Management" mean?

- **Target Scientific Applications**
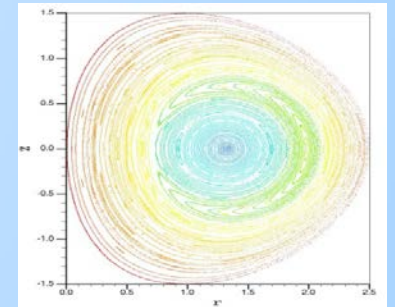  - Climate, Combustion, Fusion, Accelerator design, Cosmology,
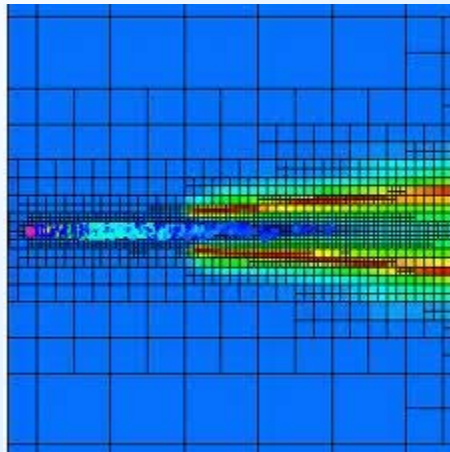




- Three pillars of science
  - Theory, Experiments, Simulations, and later Data Analysis (fourth paradigm)
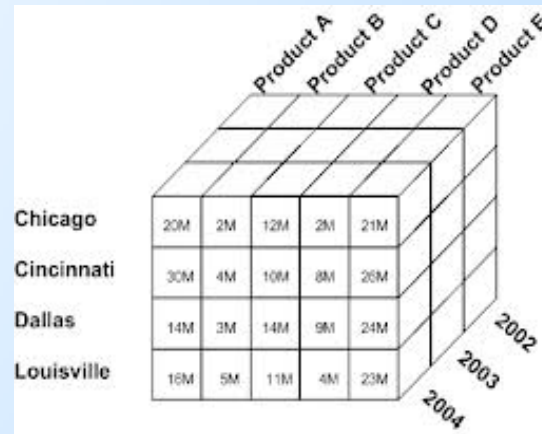


- **Algorithms, techniques, and software**
  - Representing scientific data – data models, metadata (structured/unstructured array models, geodesic models, sequence data, streaming data  )
  - Managing I/O – methods for removing I/O bottleneck
  - Accelerating efficiency of access – data structures, indexing
  - Facilitating data analysis – data manipulations for finding patterns and meaning in the data
  - Support visual analytics – accelerate extraction of subsets for real-time visualization
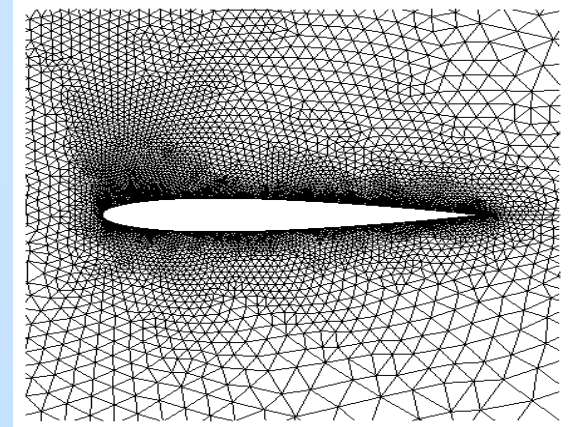
A. Shoshani
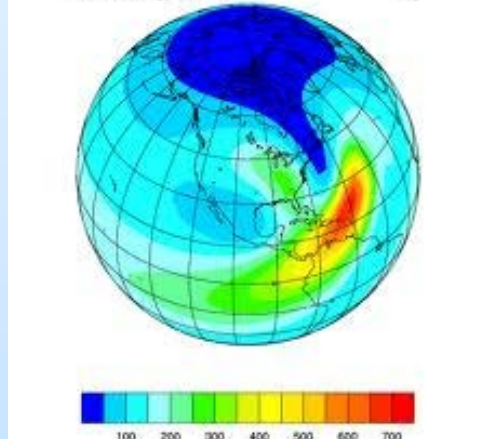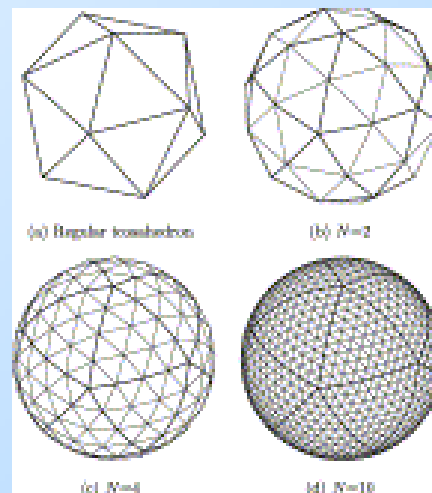
# Scientific Data Models


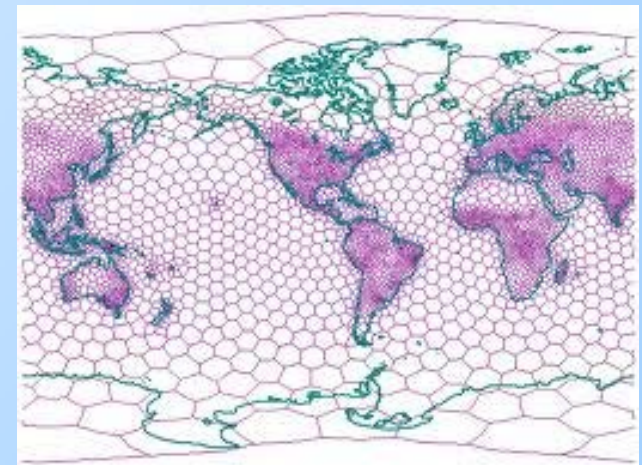Adaptive Mesh Refinement


Data Cube


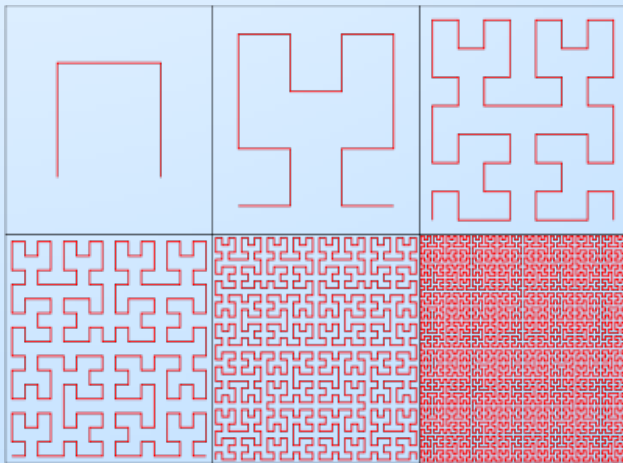Unstructured triangular grid


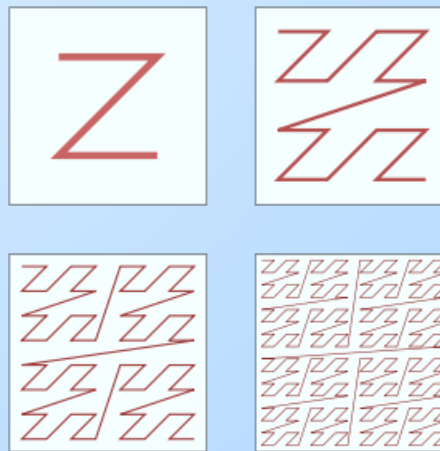Geodesic data model


Geodesic triangular data model


Unstructured grid: Voronoi tesselation

A. Shoshani

# Physical Data Structure
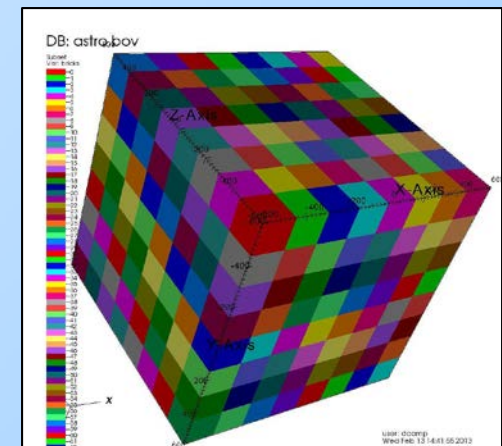
- Linearization of data based on data model
  - By coordinate order based on most prevalent access
  - Hilbert or Z-ordering to support local neighborhood access
- Partitioning data into blocks for parallel processing
  - Assigning block to different processors
  - Striping blocks on disk



Hilbert linearization order

Z-ordering

512-block dataset colored by thread ID

A. Shoshani

# Scientific data models have special operators

- Spatial structures (e.g. climate, airplane wing)
  - Region operators, slices from 3D to 2D,
- Space over time structures
  - Spatial overlap over time-steps to track pattern progress
- Temporal data
  - Before/after operators, time-overlap operators
- Time-series data (e.g. sensor data)
  - Statistical operators over regular time-intervals
- Sequence data (e.g. biology)
  - Have special alphabet (4 base-pairs for DNA, 22 for protein)
- Irregular 3D structures
  - Protein folding operators
- etc., etc.

# Scientific data management, analysis, and visualization

- **Data Management**
  -  support of physical data structures and optimization of operations over scientific logical data structures

- **Data Analysis**
  - support for manipulations of logical data structures to enhance data understanding

- **Visualization**
  - facilitating real-time visual exploration of space-time data, as well as analysis of properties of various data structures
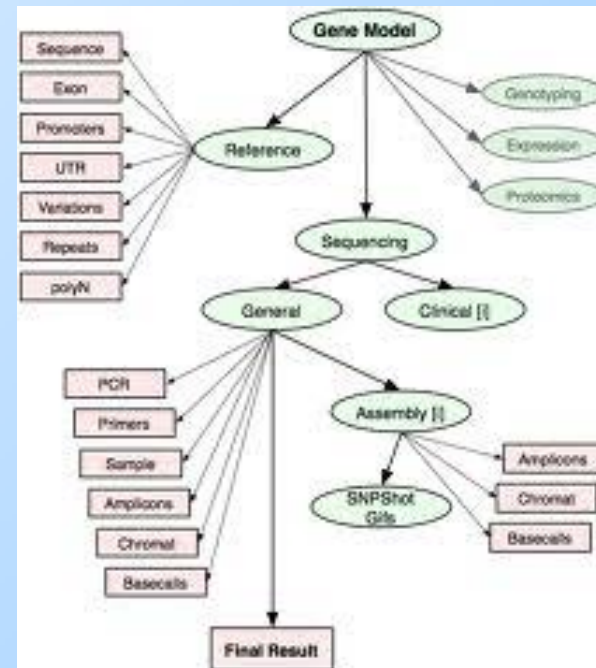
# On Scientific Metadata

Metadata is essential to describe how the data was generated/collected

- Self-describing data formats (using headers and footers) – e.g. netCDF
- Hierarchical data formats allowing organization of data as well as annotation – e.g. HDF5
- External information: who, what, when, provenance, codes, device specifics,
- Ontologies, Controlled Vocabularies



netCDF data structure



HDF5 hierarchical data format

# First SSDBM (1981) – focus on statistical data

- Menlo Park, CA

- Looking at Socio-Economic data
  - Population by (state, city, race, age, sex)
  - Socio-economic scientists did not use database systems
  - Data model does not fit relational models

- Statistical data model
  - Multi-dimensional + hierarchies over dimensions
  - Became popular with SIGMOD conferences

Statistical Data Bases Logical Model



A. Shoshani

# First SSDBM (1981) – focus on statistical data

- **OLAP**
  - Later SDBs were re-introduced as OLAP, plus operators (role-up, drill-down, )
  - Paper on "OLAP vs. Statistical Databases" – PODS 1997
  - Later OLAP was visualized as "data cubes", plus operators (Jim Gray)
  - Implementation of OLAP databases by Microsoft, Oracle, Sybase
  - Lesson: specialized systems developed for this type of a data model

- **System S**
  - 1981: Richard A. Becker:
    **Data Manipulation in the S System for Interactive Data Analysis.**
  - R is an implementation of the S programming language

LOGICAL MODEL

ROLAP REPRESENTATION

# Third SSDBM (1986) – Luxemburg

- Rojer Cubbit
- Got involved in statistical office of EU
- SSDBM started alternating between US and EU
- Introducing Scientific data
- Why? Scientists in general did not use database management systems
- VLDB 1994:
  - "Characteristics of Scientific Databases" – VLDB 1984 (Arie Shoshani, Frank Olken, Harry K. T. Wong)
  - Identified array data as an important model for scientists
  - Data kept in specialized file formats
    - NetCDF, HDF5, FITS,
    - Having their own libraries
  - This is still the case today!!!

# SSDBM (1996-1998)

- NSF got interested – Maria Zemankova
  - Suggested to alternate every year between Europe and USA
  - Before that it was every other year
- 1997 – Olympia, WA
  - Interest in Environmental Data was introduced
    Francis P. Bretherton, William L. Hibbard: Metadata: A Case Study from the Environmental Sciences.
  - Also Knowledge Discovery
    Usama M. Fayyad: Data Mining and Knowledge Discovery in Databases: Implications for Scientific Databases
  - "Summarizability" of Statistical database introduced
    Hans-Joachim Lenz, Arie Shoshani: Summarizability in OLAP and Statistical Data Bases
- 1998 – Capri
  - Interest in Multidimensional Arrays was presented
    Norbert Widmann, Peter Baumann: Efficient Execution of Operations in a DBMS for Multidimensional Arrays
  - Product: Rasdaman, open-source

# SSDBM (2001- 2004)

- 2001 – Fairfax, VA

  - Interest in Earth Systems was presented
    James Frew, Rajendra Bose: Earth System Science Workbench: A Data Management Infrastructure for Earth Science Products

- 2002 – Edinburgh

  - Interest in Biology and Gene Expression was presented
    Albert Burger, Richard A. Baldock, Yiya Yang, Andrew M. Waterhouse, Derek Houghton, Nick Burton, Duncan Davidson:
    The Edinburgh Mouse Atlas and Gene-Expression Database: A Spatio-Temporal Database for Biological Research

- 2004 – Santorini

  - Interest in Scientific Workflows was presented
    Ilkay Altintas, Chad Berkley, Efrat Jaeger, Matthew B. Jones, Bertram Ludäscher, Steve Mock: Kepler: An Extensible System for Design and Execution of Scientific Workflows

  - Led to Kepler, an open-source product

A. Shoshani

# SSDBM (2008- 2012)

- 2008 – Hong Kong
  - Interest in Scientific Ontology Databases was presented
    Paea LePendu, Dejing Dou, Gwen A. Frishkoff, Jiawei Rong: Ontology Database: A New Method for Semantic Modeling and an Application to Brainwave Data.

- 2011 – Portland
  - Interest in Scientific Database Systems was presented
    Michael Stonebraker, Paul Brown, Alex Poliakov, Suchi Raman: The Architecture of SciDB
  - Product: open source SciDB

- 2012 – Crete
  - Interest in Data Fusion was presented
    David Maier, V. M. Megler, António M. Baptista, Alex Jaramillo, Charles Seaton, Paul J. Turner: Navigating Oceans of Data.

# SSDBM (2013- 2017)

- 2013 – Baltimore

  - Interest in Big Data was presented (keynote)
    Michael J. Franklin:  Making Sense of Big Data with the Berkeley Data Analytics Stack.

  - Interest in Streaming Data was presented
    Hamid Mousavi*, Carlo Zaniolo:  Fast Computation of Approximate Biased Histograms on Sliding Windows over Data Streams

- 2016 – Budapest

  - Interest in User-Defined-Functions (UDF) was presented
    Mark Raasveldt, Hannes Mühleisen: Vectorized UDFs in Column

  - Stores

  - Product: open source MonetDB/Python

- 2017 – Chicago

  - Interest in N-dimensional Arrays was presented
    Veranika Liaukevich, Dimitar Mišev, Peter Baumann, Vlad Merticariu: Location and Processing Aware Datacube Caching

# Final Thoughts

- Work with domain scientists and identify their data problems
  - Their logical/abstract data model
  - Their operators on that data models, including functions on the data (UDFs)
  - Their metadata, ontology, controlled vocabularies
  - Their data constraints
- Finds out how they store their data – specialized file formats
  - Do not try to force them to reshaped their data into your system (too big of a task, they will loose interest)
- Build something useful to them, and integrate in their environment
  - That will keep their attention for continued collaboration
- Submit your paper(s) to SSDBM 🙂

# LIGHT-HEARTED ANECDOTES

A. Shoshani

# Fun memories

- 1988, Rome, Italy
  - Gucci bags to all + Channel perfume for woman
  - Banquet at an estate outside Rome, six course
- 1996, Stockholm, Sweden
  - River ride to forest, walk to banquet
- 1997, Olympia, Washington
  - Nature walk to ocean
- 1998, Capri, Italy
  - One afternoon free to visit blue grotto
- 2000, Berlin, Germany
  - River boat ride
- 2002, Edinburgh, Scotland
  - Yearly fireworks spectacular display
- 2004, Santorini, Greece
  - Boat ride to the islands – swimming in sea
- 2005, Santa Barbara, California
  - Banquet: barbecue on beach
- 2007, Banff, Canada
  - Spectacular nature setting in the park

- 2010, Heidelberg, Germany
  - Held at European Media Lab – beautiful gardens
- 2011, Portland, Oregon
  - Great beer at location near river
- 2012, Crete, Greece
  - Beautiful hotel with view of Mediterranean

A. Shoshani
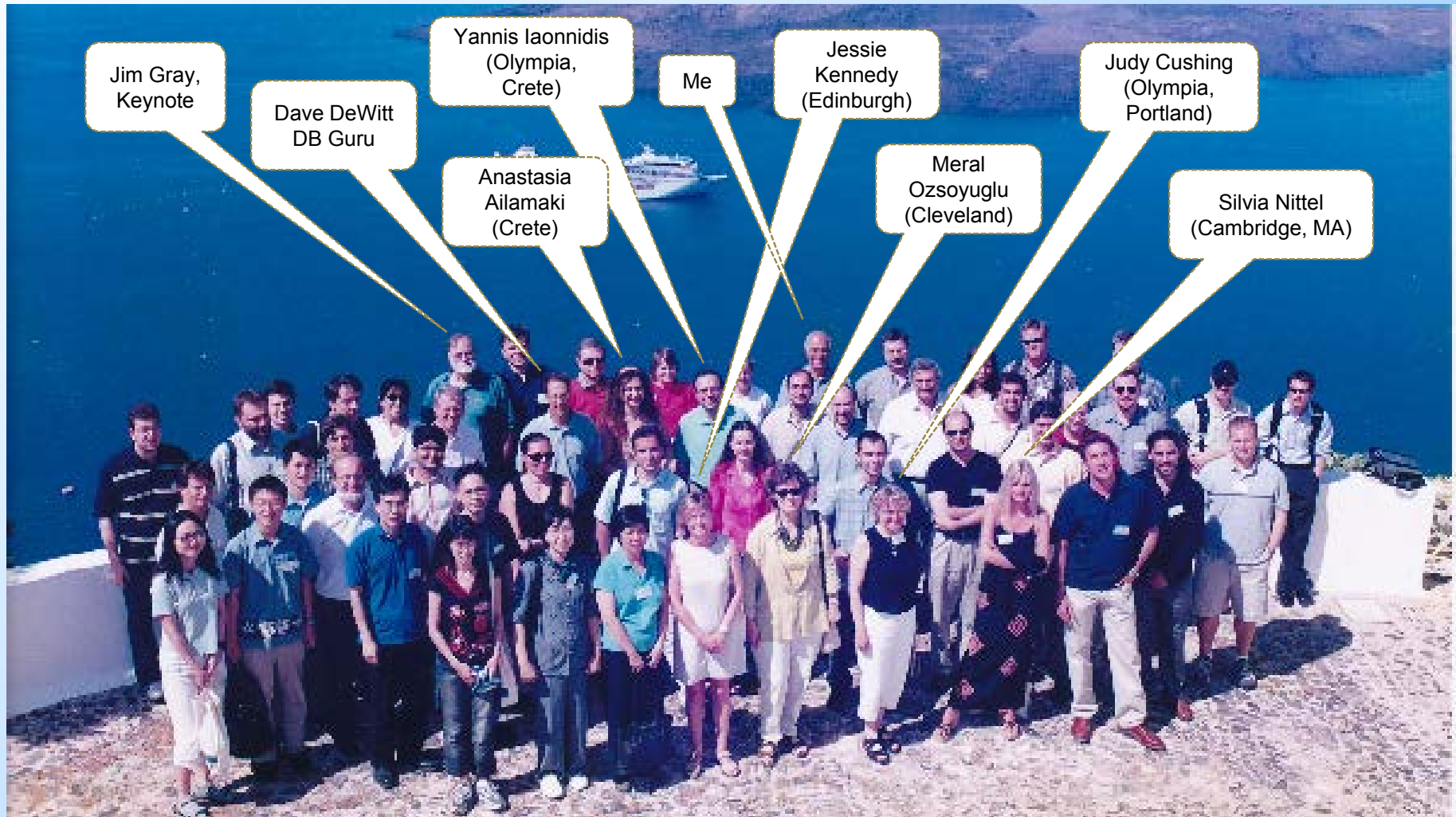
# Berlin (2000)



Yannis Iaonnidis (Olympia, Crete)

Me

Jessie Kennedy (Edinburgh)

A. Shoshani

# Santorini (2004)



A. Shoshani

# Banff (2007)



Marianne
Winslett
(New Orleans)

# New Orleans (2009)



21st International Conference
on
Scientific and Statistical
Database Management

June 2–4, 2009

Room: Vieux Carré
17th Floor

A. Shoshani

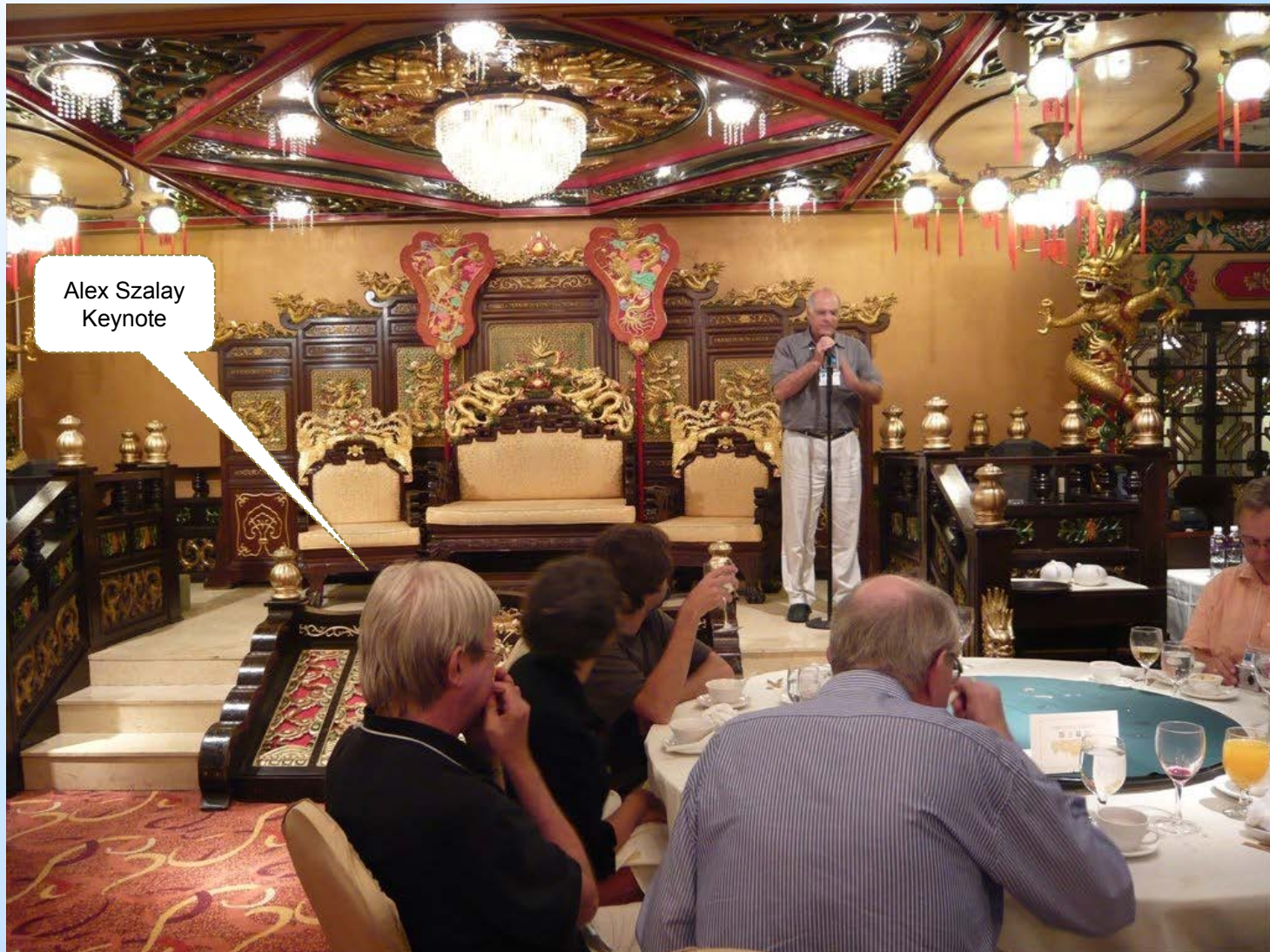# Crete (2012)



A. Shoshani

# Hong Kong (2008)

Crown Chair

# Hong Kong (2008)



A. Shoshani

# Hong Kong – rain, rain, everywhere



Alex Szalay
Keynote

# Aalborg (2014)



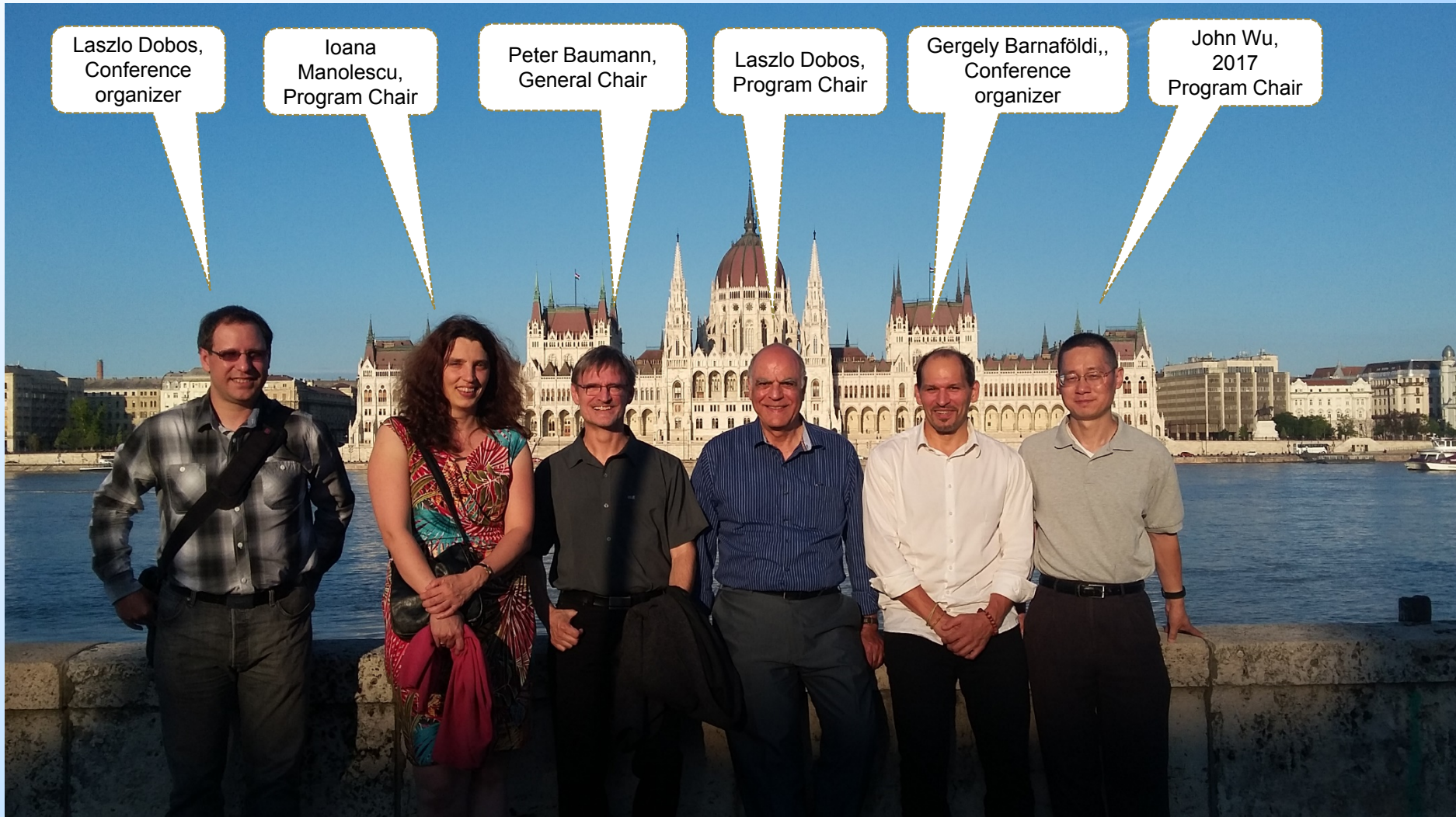Torben Pedersen, Program Chair

# Budapest (2016)



A. Shoshani

# THE END