



Network Measurement Virtual Observatory: An Integrated Database Environment for Internet Measurements and Data Analysis

Tamás Sebök, Zsófia Kallus, Sándor Laki,
Péter Mátray, József Stéger, Péter Hága,
László Dobos, István Csabai, Gábor Vattay

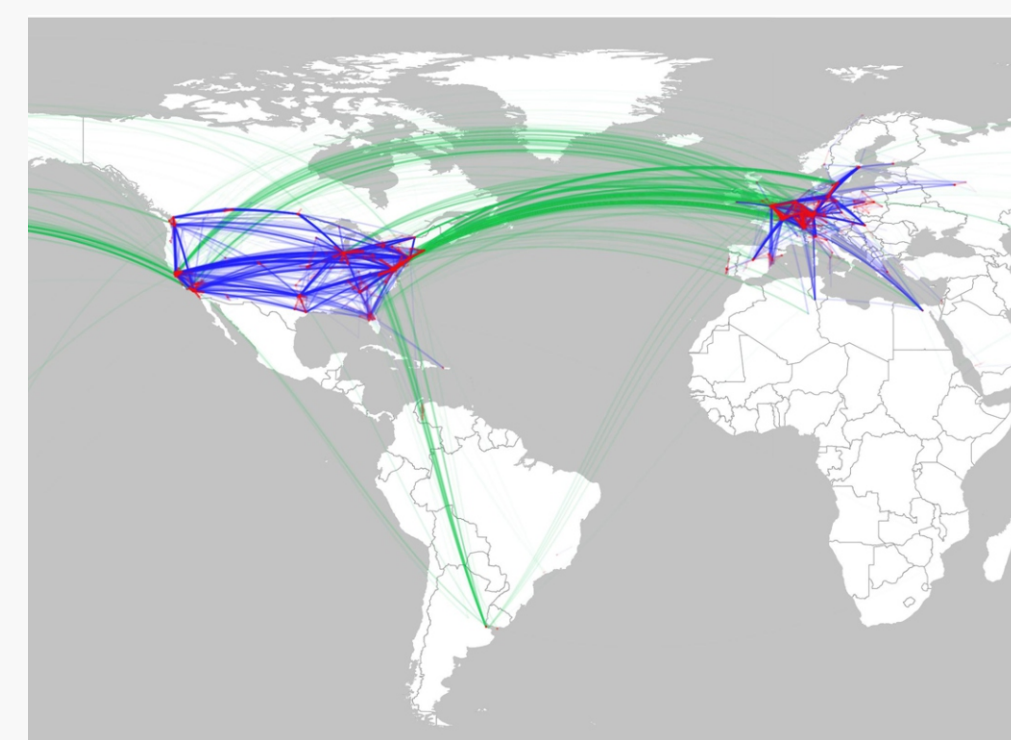
Eötvös Loránd University, Department of Physics of Complex Systems, Budapest, Hungary

Internet tomography techniques work by setting up servers with precisely synchronized clocks all around the world to measure ping delay times and to determine the structure and dynamics of the multiply connected computer network. The experiments require coordination of hundreds of computers, billions of measurements and software capable of gathering and analyzing measurement data. The high number of these micro-experiments makes Internet tomography a particularly interesting big data science.

We present the Network Measurement Virtual Observatory (nmVO), an integrated system built around a relational database system, custom web services and various open software developed primarily for other scientific fields. Its aim is to serve collaborating research groups sharing their measurement results through standardized archives. The tight integration of the data acquisition process with data analysis and the final data products makes it possible to perform historic and real-time observations of the Internet using simple database queries. A case study of the nmVO is our geographic localization service called Spotter.

X-raying the Internet

- Internet tomography
 - Build graph structure from traceroutes and ping delay data
 - Topology of the network
 - Bandwidth along edges



The Concept of nmVO

- Virtual Observatories: data warehouses for scientific big data
- nmVO: a dynamic data warehouse:
 - network measurements: millions of micro-measurements
 - conducting measurement requires large amount of input data
 - raw data collection directly into the database
 - data reduction is done inside the database
 - new data gets published instantly after reduction
 - re-analyze measurements later, if necessary
 - use archival data to analyze long-term dynamics of Internet
- nmVO is built on well-known scientific database software
 - CasJobs, originally built for the astronomical SkyServer
 - primary user interface: SQL-based batch job system
 - Hierarchical Triangular Mesh used for geographical data

nmVO in practice

- Multi-TB SQL database
 - collect and store raw data
- Reduction and analysis tasks are written in SQL
 - run data reduction and analysis inside DB
 - store results in the same DB
 - users can run the processing on the server where data reside
- Integrated with data collection services
 - initiate new Internet measurements right from SQL!
 - great potential, e.g. collect missing data on-the-fly
- Access to remote datasets helps incorporate data from other Internet measurement projects

nmVO infrastructure

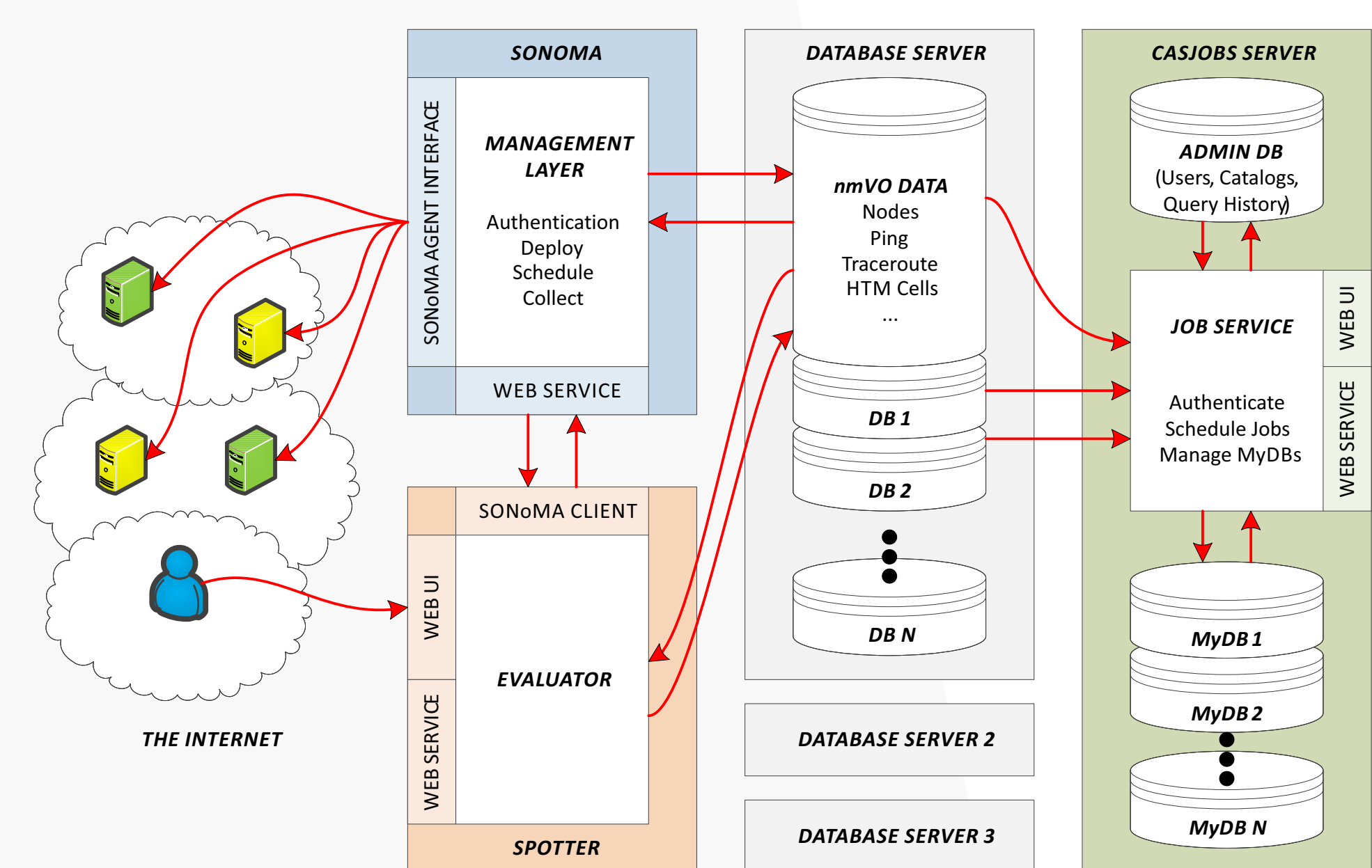


Figure: The nmVO architecture. Services are centered around the database servers (grey) storing raw and processed measurement data. The SONOMA management layer (blue) is responsible for conducting measurements by controlling measurement software on remote server nodes around the world. The Spotter web service (orange) is built on top of SONOMA and nmVO. When a user initiates an IP localization, Spotter instructs SONOMA to collect necessary ping delay data and analyzes the raw measurements. Both raw data and analyzed measurements are made instantly public in nmVO. Users can interact with the system via CasJobs (green), a SQL query interface and batch execution layer. Data processing tasks are formulated in SQL and executed by the database servers. SONOMA functions are directly available as SQL stored procedures which allow on-the-fly experiments.

SONOMA

- Use Etomic and PlanetLab servers around the world
- SOAP-based web services to collect Internet measurement data
 - hides complexities of underlying infrastructure
 - easy to implement interfaces for network experiment designers
- Wide set of atomic measurements
 - ping, tracerout, chirp, train, DNS lookup etc.

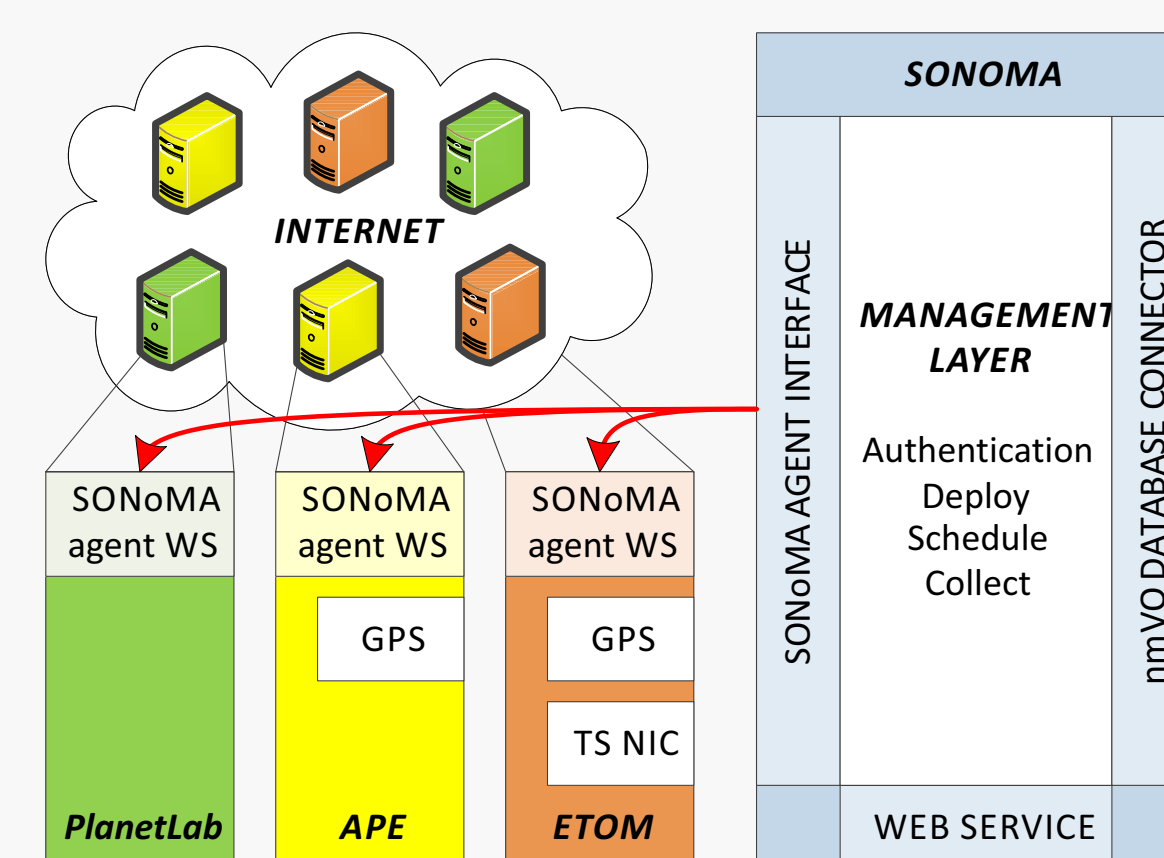


Figure: The SONOMA architecture. SONOMA agents run on network measurement servers world wide. SONOMA currently supports the PlanetLab and Etomic infrastructures. Etomic machines are also equipped with high precision GPS clocks for accurate timing of ping delays. The SONOMA management layer (grey) is responsible for conducting measurements and store raw and processed data directly in nmVO. SONOMA is accessible via its SOAP web service interface.

Spotter

- Geolocation service built on nmVO and SONOMA
- Determines location of host from its IP address
 - landmarks: servers with known coordinates
 - ping delay times measured from landmarks
 - location determined from ping delay times
- Raw data and measurements published in nmVO

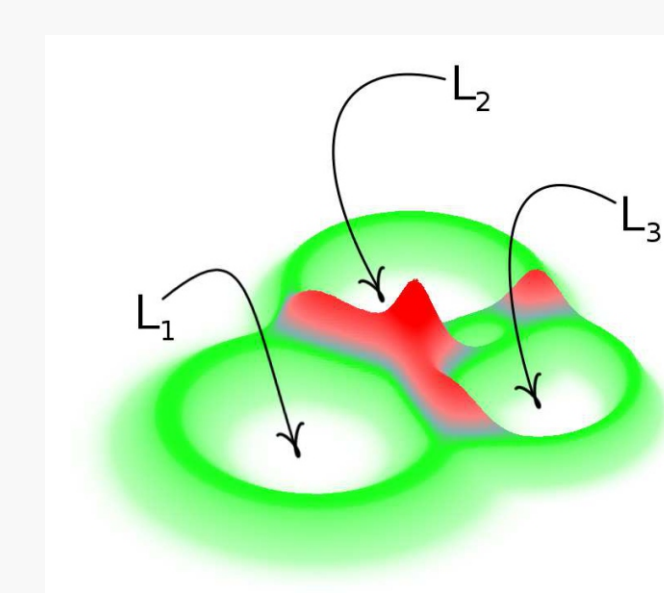
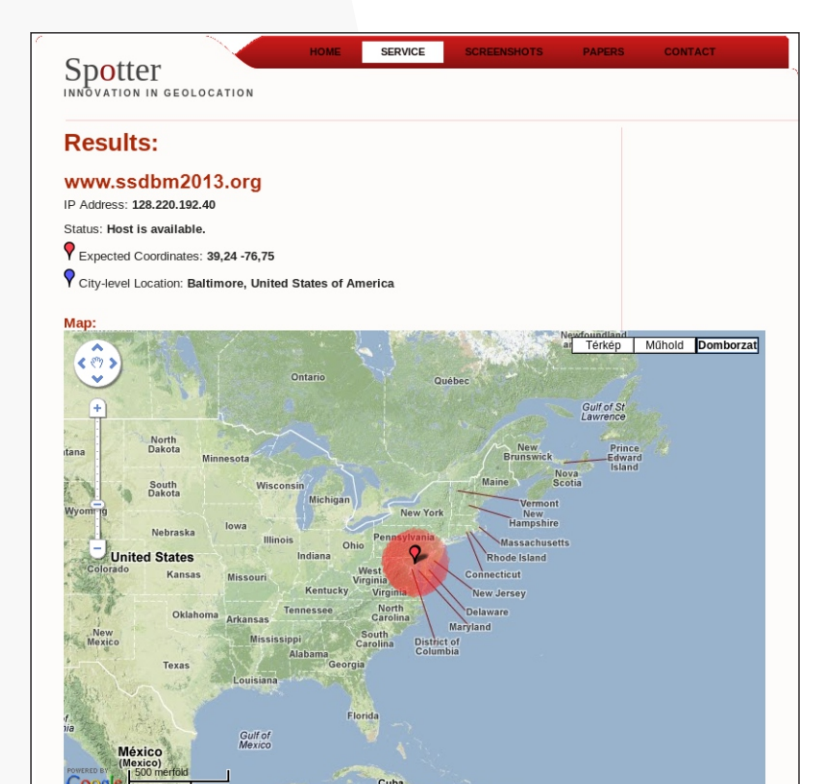


Figure: Results page of Spotter localizing ssdbm2013.org. The accuracy of geolocation depends significantly on the available landmarks (servers with known coordinates).

Figure: Spotter uses a complex statistical model based on geographic coordinates of landmarks and the measured ping delay times is used to calculate the PDF of the approximate coordinates of the localized IP address.



Please visit <http://nm.vo.elte.hu>