

Wright State University

## CORE Scholar

---

Computer Science and Engineering Faculty  
Publications

Computer Science & Engineering

---

2016

### HPC Enabled Data Analytics for High-Throughput High-Content Cellular Analysis

Ross A. Smith

Rhonda J. Vickery

Jack Harris

Sara Gharabaghi

*Wright State University - Main Campus, gharabaghi.2@wright.edu*

Thomas Wischgoll

*Wright State University - Main Campus, thomas.wischgoll@wright.edu*

*See next page for additional authors*

Follow this and additional works at: <https://corescholar.libraries.wright.edu/cse>



Part of the [Computer Sciences Commons](#), and the [Engineering Commons](#)

---

#### Repository Citation

Smith, R. A., Vickery, R. J., Harris, J., Gharabaghi, S., Wischgoll, T., Short, D., Trevino, R., Kawamoto, S. A., Lamkin, T. J., Schoen, K., Bardes, E. E., Tabar, S. C., & Aronow, B. J. (2016). HPC Enabled Data Analytics for High-Throughput High-Content Cellular Analysis. .  
<https://corescholar.libraries.wright.edu/cse/469>

This Conference Proceeding is brought to you for free and open access by Wright State University's CORE Scholar. It has been accepted for inclusion in Computer Science and Engineering Faculty Publications by an authorized administrator of CORE Scholar. For more information, please contact [library-corescholar@wright.edu](mailto:library-corescholar@wright.edu).

---

**Authors**

Ross A. Smith, Rhonda J. Vickery, Jack Harris, Sara Gharabaghi, Thomas Wischgoll, David Short, Robert Trevino, Steven A. Kawamoto, Thomas J. Lamkin, Kevin Schoen, Eric E. Bardes, Scott C. Tabar, and Bruce J. Aronow

# HPC Enabled Data Analytics for High-Throughput High-Content Cellular Analysis

Ross A. Smith, Rhonda J. Vickery, Jack Harris  
Engility Corporation  
Dayton, Ohio, U.S.A.  
Ross.Smith@engilitycorp.com

Robert Trevino, Steven A. Kawamoto,  
Thomas J. Lamkin, Kevin Schoen  
U.S. Air Force Research Laboratory  
WPAFB, Ohio, U.S.A.

Sara Gharabaghi, Thomas Wischgoll, David Short  
Wright State University  
Dayton, Ohio, U.S.A.

Eric E. Bardes, Scott C. Tabar, Bruce J. Aronow  
Cincinnati Children's Hospital Medical Center  
Cincinnati, Ohio, U.S.A.

**Abstract**—Biologists doing high-throughput high-content cellular analysis are generally not computer scientists or high performance computing (HPC) experts, and they want their workflow to support their science without having to be. We describe a new HPC enabled data analytics workflow with a web interface, HPC pipeline for analysis, and both traditional and new analytics tools to help them transition from a single workstation mode of operation to power HPC users. This allows the processing of multiple plates over a short period of time to ensure timely query and analysis to match potential countermeasures to individual responses.

**Keywords**—High-Throughput High-Content Screening; Pipeline; Analytics; Visualization.

The INSIGHTS project discovers beneficial chemical or genetic interrogations for human performance, medical intelligence, or therapeutic application [1]. The envisioned use cases are force protection of U.S. warfighters from natural exposure or biological attack. Rapid screening of numerous possible countermeasures decreases the time-to-decision from years to months or even weeks. The methodology is based on high-throughput, high-content biological screening and matching potential countermeasures to individual response. Individual experiments are done within wells of an experimental plate (see Fig. 1). After incubation, each well is observed at multiple sites, and resulting images are submitted to the pipeline for cell segmentation (CS), generating image masks to identify each cell in a site. Each cell then goes through featurization, generating 11,000+ features. Feature selection removes all but the most informative features. Finally, well scoring identifies interrogations for further investigation.

Each experimental batch has multiple plates where each well is imaged at multiple sites at frequencies optimized with commonly used biological fluorescent stains or when using non-fluorescent phase contrast imaging. This approach typically yields four images per site and 3,000 sites produce 12,000 images per plate. Although sites have around 50 cells, more than 300 cells are not uncommon. Processing 11,000 features for

hundreds of plates, including additional analytics, requires HPC and produces millions of files. After the images have been collected, they are uploaded to a file system mounted by the HPC head nodes as shown in Fig. 2. The researcher ensures batch metadata is linked via a project database and launches analysis jobs via an easy to use  $\mu$ Batch web portal, where clients on HPC nodes retrieve  $\mu$ Batch jobs following a “bag of tasks” paradigm. The clients written in C++, Python, Java, and C# execute the data analytics pipeline, updating both the project database and the  $\mu$ Batch server.

The biologist has additional flexibility to do post analysis to facilitate verification and validation of results. An example is the ability to create traditional dose response curve plots from downloaded Kolmogorov–Smirnov reports, as shown in Fig. 3 with the first 25 of 32 plots displayed for a plate. Dimple, based

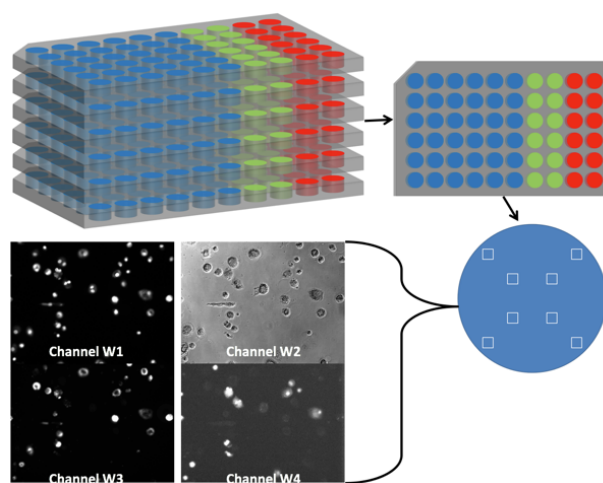


Fig. 1. An experimental batch can have multiple plates. Each plate has 384 wells. Each well is sampled at eight sites. Four images are captured at each site based on filters and dye frequencies.

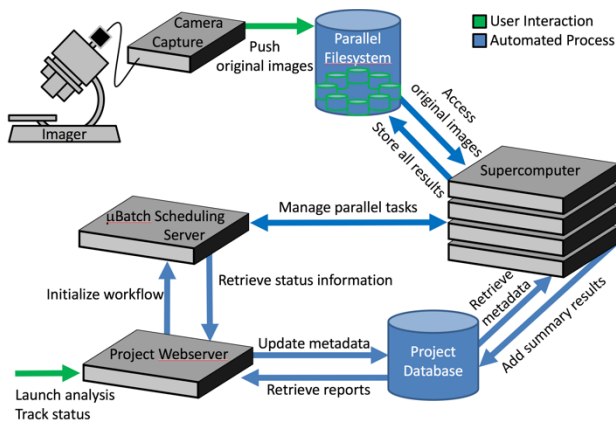


Fig. 2. Data movement through the system, including files, database accesses, job launches, and user input. The 12,000+ images from a single plate take up 33 GB. Compressed feature data takes up 3 GB or more depending on the number of cells.

on Data Driven Documents (D3), was used for the visualization, and NW.js was used to prototype the web page as a desktop application.

The blue circles indicate actual points from the data and hovering over one of these shows the data values for that point, as shown in the center plot. The red curves indicate where dose response curves were found using the Hill equation, and a red dashed line marks the half maximal inhibitory concentration (IC50). This is a common analysis for comparison of individual treatments or responses. This is supplemented with a D3 visualization showing a more comprehensive comparison of features by logical grouping using parallel coordinates plots

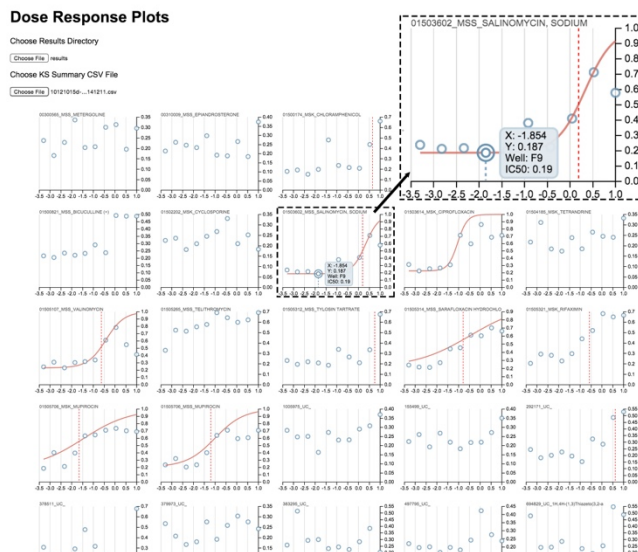


Fig. 3. Dose response curves with single point highlighted.

with controls to determine which treatments cause responses that are dissimilar to uninfected cells, as shown in Fig. 4 [2]. In this way, both familiar and more comprehensive tools are available.

We are currently adapting the capability to other biology and chemistry oriented problem datasets. A public facing portal is also being developed.

#### ACKNOWLEDGMENTS

The authors would like to acknowledge the computational resources and PETTT software support from the DoD High Performance Computing Modernization office under Contract No. GS04T09DBC0017 and biological support for molecular signatures from the U.S. Air Force under Contract No. FA8650-14D-6516.

#### REFERENCES

- [1] Vickery, R.J., Smith, R., Bardes, E., Tabar, S., Short, D., and Trevino, R., Utilization of Hybrid Computing for High Throughput Identification of Beneficial Chemicals or Biological Interrogations for Human Effectiveness, PETTT PP-ACE-KY06-003-P3 Final Technical Report, 21 August 2015.
- [2] Wischgoll, T., Vickery, R., and Smith, R., Extending the Visual Analytics Framework for High Throughput Screening of Biological Infectious Agents Project, PETTT PP-ACE-KY05-005-P3 Final Technical Report, 22 August 2014.

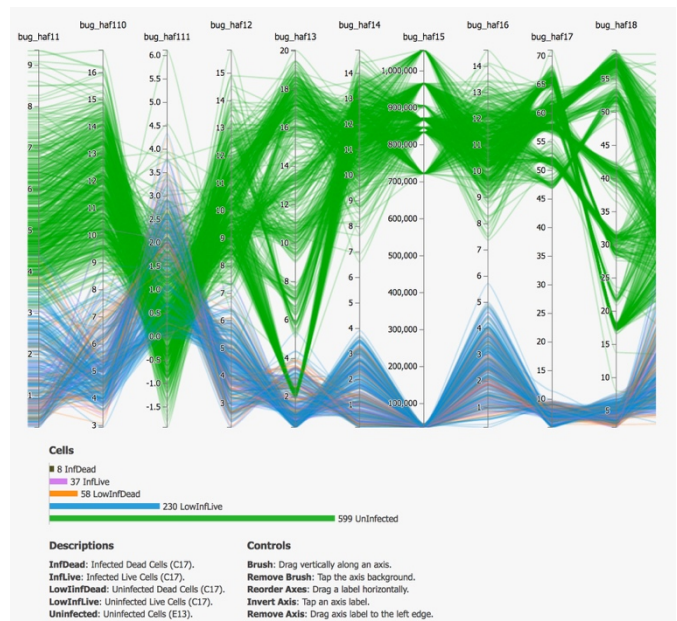


Fig. 4. Parallel coordinate visualization of Haralick features extracted from cell images.