# Real-World Analytics



**David A. Bader**

🐦 **@Prof_DavidBader**

http://www.cs.njit.edu/~bader



New Jersey Institute
of Technology

# David A. Bader

## Distinguished Professor and
## Director, Institute for Data Science

- IEEE Fellow, SIAM Fellow, AAAS Fellow

- Recent Service:
  - White House's National Strategic Computing Initiative (NSCI) panel
  - Computing Research Association Board
  - NSF Advisory Committee on Cyberinfrastructure
  - Council on Competitiveness HPC Advisory Committee
  - IEEE Computer Society Board of Governors
  - IEEE IPDPS Steering Committee
  - Editor-in-Chief, ACM Transactions on Parallel Computing
  - Editor-in-Chief, IEEE Transactions on Parallel and Distributed Systems

- Over $183M of research awards

- 230+ publications, ≥ 8,400 citations, h-index ≥ 54

- National Science Foundation CAREER Award recipient

- Directed: NVIDIA GPU Center of Excellence

- Directed: Sony-Toshiba-IBM Center for the Cell/B.E. Processor

- Founder: Graph500 List benchmarking "Big Data" platforms

- Recognized as a "RockStar" of High Performance Computing by InsideHPC in 2012 and as HPCwire's People to Watch in 2012 and 2014.

NJIT
New Jersey Institute
of Technology

# NJIT Jumps into Top 100 for 2020 U.S. News College Rankings



http://news.njit.edu/njit-jumps-top-100-2020-us-news-college-rankings

9 September 2019

## America's Great Working-Class Colleges

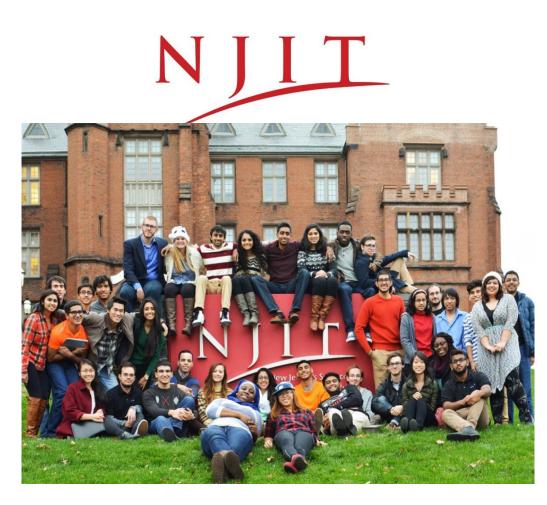David Leonhardt

# An Upward Mobility Top 10

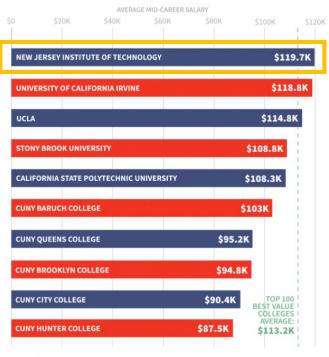Colleges ranked by percent of students from the bottom fifth of the income scale who end up in the top three-fifths.

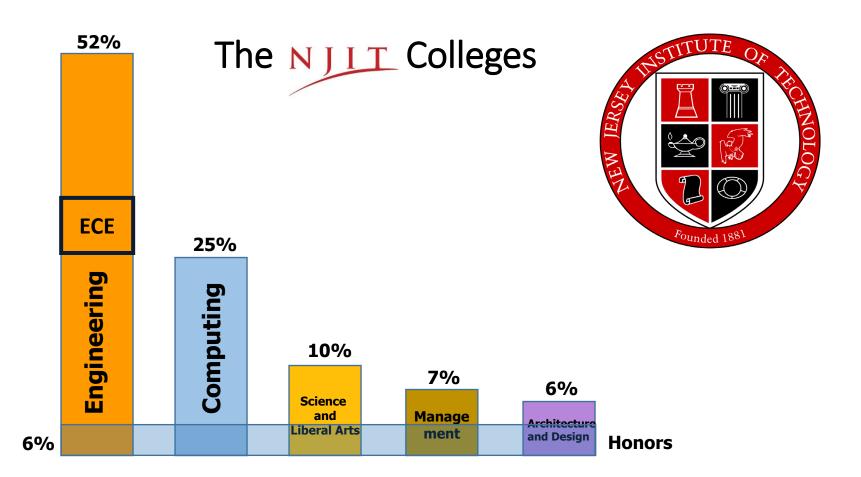| 1 | New Jersey Institute of Technology | 85% |
|---|---|---|
| 2 | Pace University | 82% |
| 3 | Calif. State – Bakersfield | 82% |
| 4 | Univ. California – Irvine | 81% |
| 5 | Calif. Poly – Pomona | 81% |
| 6 | Xavier of Louisiana | 80% |
| 7 | SUNY Stony Brook | 79% |
| 8 | San Jose State | 79% |
| 9 | CUNY Baruch College | 79% |
| 10 | Calif. State – Long Beach | 78% |

# NJIT

# Forbes

April 10, 2018

## AMERICAN DREAM U
### BEST VALUE COLLEGES WITH THE HIGHEST UPWARD MOBILITY RATES

We looked at the top 100 Forbes Best Value Colleges 2018 to identify the 10 schools that had the most success with upward mobility—the percentage of students from the bottom 20% income distribution who reach the top 20%—and then ranked them by alumni with the highest mid-career salaries.

AVERAGE MID-CAREER SALARY

| School | Salary |
|---|---|
| NEW JERSEY INSTITUTE OF TECHNOLOGY | $119.7K |
| UNIVERSITY OF CALIFORNIA IRVINE | $118.8K |
| UCLA | $114.8K |
| STONY BROOK UNIVERSITY | $108.8K |
| CALIFORNIA STATE POLYTECHNIC UNIVERSITY | $108.3K |
| CUNY BARUCH COLLEGE | $103K |
| CUNY QUEENS COLLEGE | $95.2K |
| CUNY BROOKLYN COLLEGE | $94.8K |
| CUNY CITY COLLEGE | $90.4K |
| CUNY HUNTER COLLEGE | $87.5K |

TOP 100 BEST VALUE COLLEGES AVERAGE: $113.2K

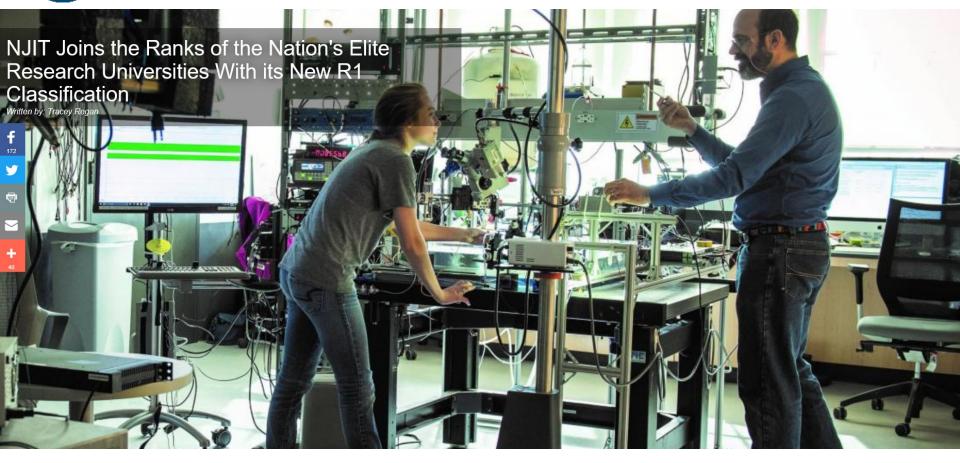DATA: The Equal Opportunity Project, PayScale

The NJIT Colleges

# Computing at NJIT



- Largest computing program in the NY/NJ metro area

- All undergraduate students take a capstone project course – working with companies

- Active internships and co-op programs in cooperation with industry, government and others

- Student clubs organize extra-curricula activities, e.g. hackathons and field trips, and provide tutoring

8

# Industrial Partners

# The Carnegie Classification of Institutions of Higher Education ®

NJIT Joins the Ranks of the Nation's Elite Research Universities With its New R1 Classification

Written by: Tracey Regan

f
172

y

🖶

✉

+
40

# Cyber Innovations for Solving Global Grand Challenges



New Jersey Institute of Technology



LANL Roadrunner with IBM Cell B.E. Top500 No. 1 system from June 2008 to June 2009



Intel HIVE processor (2019)



IBM Watson with POWER7/8, won Jeopardy in Feb 2010



Echelon: NVIDIA & Team's UHPC Project
STEVE KECKLER
DIRECTOR OF ARCHITECTURE RESEARCH, NVIDIA

NVIDIA GPUs used in 127 of Top500 systems, incl. top 2 (in USA), and fastest in Europe and Japan. (Nov. 2018)



Cray XMT with ThreadStorm proc. Massively Multithreaded Architecture



IBM BlueGene/Q. Record breaking performance over 10PF sustained on science apps

# Data Science: Discovery and Innovation

The ability to manipulate data and understand Data Science is becoming increasingly critical to current and future discovery and innovation.
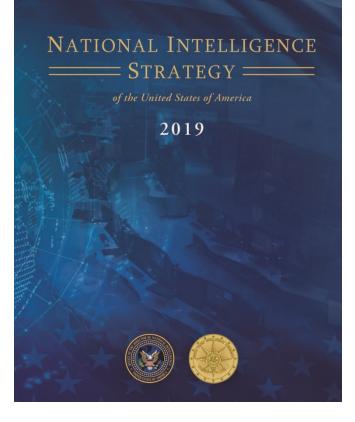
REALIZING THE POTENTIAL OF DATA SCIENCE Final Report from the National Science Foundation Computer and Information Science and Engineering Advisory Committee Data Science Working Group. Francine Berman and Rob Rutenbar, co-Chairs Henrik Christensen, Susan Davidson, Deborah Estrin, Michael Franklin, Brent Hailpern, Margaret Martonosi, Padma Raghavan, Victoria Stodden, Alex Szalay. December 2016

**The National Strategic Computing Initiative NSCI**

The National Strategic Computing Initiative (NSCI) The NSCI was launched by Executive Order (EO) 13702 in July 2015 to advance U.S. leadership in high performance computing (HPC).

McKinsey predicts that data-driven technologies will bring an additional $300 billion of value to the U.S. health care sector alone, and by 2020, 1.5 million more "data-savvy managers" will be needed to capitalize on the potential of data, "big" and otherwise.

Manyika, J. et al. (2011). Big data: The next frontier for innovation, competition, and productivity. McKinsey Global Institute. Retrieved from http://www.mckinsey.com/business-functions/business-technology/our-insights/big-data-the-next-frontier-for-innovation
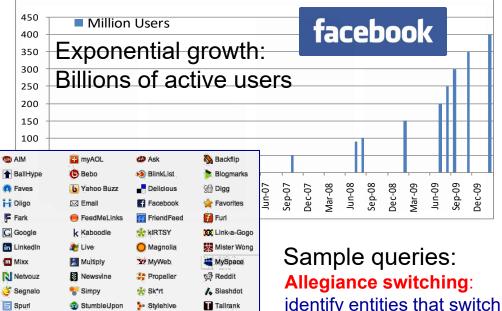
# Strategic Intelligence



"Advances in communications and the democratization of other technologies have also generated an ability to create and share vast and exponentially growing amounts of information farther and faster than ever before. This **abundance of data** provides significant opportunities for the IC, including new avenues for collection and the potential for greater insight, but it also challenges the IC's ability to collect, process, evaluate, and analyze such enormous volumes of data quickly enough to provide relevant and useful insight to its customers."

→ "Develop and maintain capabilities to acquire and evaluate data to obtain a deep understanding of the global political, diplomatic, military, economic, security, and informational environment. "

# Exascale Streaming Data Analytics: Real-world challenges

**All involve analyzing massive streaming complex networks:**

- **Health care** → disease spread, detection and prevention of epidemics/pandemics (e.g. SARS, Avian flu, H1N1 "swine" flu)

- **Massive social networks** → understanding communities, intentions, population dynamics, pandemic spread, transportation and evacuation

- **Intelligence** → business analytics, anomaly detection, security, knowledge discovery from massive data sets

- **Systems Biology** → understanding complex life systems, drug design, microbial research, unravel the mysteries of the HIV virus; understand life, disease,

- **Electric Power Grid** → communication, transportation, energy, water, food supply

- **Modeling and Simulation** → Perform full-scale economic-social-political simulations



Exponential growth: Billions of active users

Ex: discovered minimal changes in O(billions)-size complex network that could hide or reveal top influencers in the community

Sample queries:
**Allegiance switching**: identify entities that switch communities.
**Community structure**: identify the genesis and dissipation of communities
**Phase change**: identify significant change in the network structure

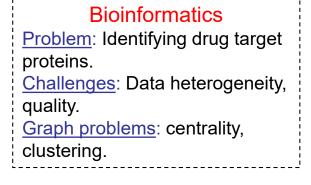**REQUIRES PREDICTING / INFLUENCE CHANGE IN REAL-TIME AT SCALE**

# Graphs are pervasive in large-scale data analysis

- **Sources** of massive data: peta- and exa-scale simulations, experimental devices, the Internet, scientific applications.

- **New challenges for analysis**: data sizes, heterogeneity, uncertainty, data quality.

**Astrophysics**
Problem: Outlier detection.
Challenges: massive datasets, temporal variations.
Graph problems: clustering, matching.

**Bioinformatics**
Problem: Identifying drug target proteins.
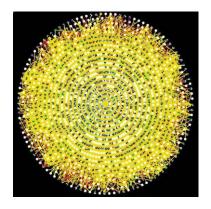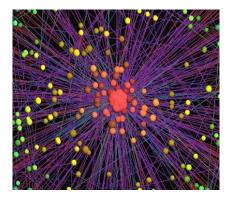Challenges: Data heterogeneity, quality.
Graph problems: centrality, clustering.

**Social Informatics**
Problem: Discover emergent communities, model spread of information.
Challenges: new analytics routines, uncertainty in data.
Graph problems: clustering, shortest paths, flows.





Image sources: (1) http://physics.nmt.edu/images/astro/hst_starfield.jpg
(2,3) www.visualComplexity.com

# Network Analysis for Intelligence and Survelliance

- [Krebs '04] Post 9/11 Terrorist Network Analysis from public domain information

- Plot masterminds correctly identified from interaction patterns: centrality



Image Source: http://www.orgnet.com/hijackers.html

- A global view of entities is often more insightful

- Detect anomalous activities by exact/approximate graph matching



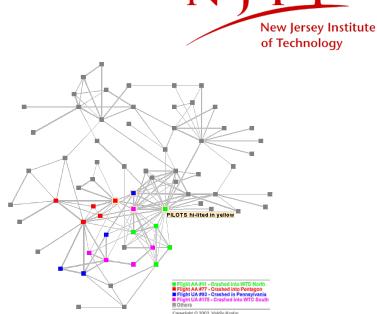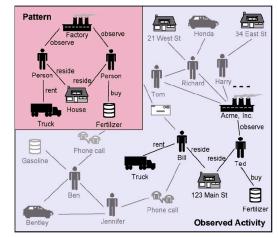Image Source: T. Coffman, S. Greenblatt, S. Marcus, Graph-based technologies for intelligence analysis, CACM, 47 (3, March 2004): pp 45-47

# Characterizing Graph-theoretic computations

**Input: Graph abstraction**

**Problem: Find \*\*\***

- paths
- clusters
- partitions
- matchings
- patterns
- orderings

**Graph algorithms**

- traversal
- shortest path algorithms
- flow algorithms
- spanning tree algorithms
- topological sort
  .....

**Factors that influence choice of algorithm**

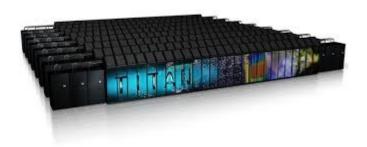- graph sparsity (m/n ratio)
- static/dynamic nature
- weighted/unweighted, weight distribution
- vertex degree distribution
- directed/undirected
- simple/multi/hyper graph
- problem size
- granularity of computation at nodes/edges
- domain-specific characteristics

Graph problems are often recast as sparse linear algebra (e.g., partitioning) or linear programming (e.g., matching) computations

# Streaming Analytics move us from reporting the news to predictive analytics

## Traditional HPC

- Great for "static" data sets.

- Massive scalability at the cost of programmability.

- Great for dense problems.

  - Sparse problems typically underutilize the system.

## Streaming Analytics

- Requires specialized analytics and data structures.

- Rapidly changing data.

- Low data re-usage.

  - Focused on memory operations and not FLOPS.

# Massive Streaming Graph Analytics



Analysts

(A, B, t1, poke)
(A, C, t2, msg)
(A, D, t3, view wall)
(A, D, t4, post)

(B, A, t2, poke)
(B, A, t3, view wall)
(B, A, t4, msg)

Q3? Q2? Q1?

... e9   e8   e7   e6   e5   e4   e3   e2   e1 ...
Billions of edges

David A. Bader

# Mining Twitter for Social Good

ICPP 2010

## Massive Social Network Analysis: Mining Twitter for Social Good
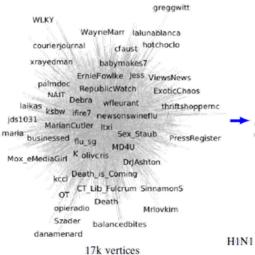
David Ediger   Karl Jiang
Jason Riedy   David A. Bader
Georgia Institute of Technology
Atlanta, GA, USA

Courtney Corley   Rob Farber
Pacific Northwest National Lab.
Richland, WA, USA

William N. Reynolds
Least Squares Software, In
Albuquerque, NM, USA

*Abstract*—Social networks produce an enormous quantity of data. Facebook consists of over 400 million active users sharing over 5 *billion* pieces of information each month. Analyzing this vast quantity of unstructured data presents challenges for software and hardware. We present GraphCT, a *Graph Characterization Toolkit* for massive graphs representing social network data. On a 128-processor Cray XMT, GraphCT estimates the betweenness centrality of an artificially generated (R-MAT) 537 million vertex, 8.6 billion edge graph in 55 minutes and a real-world graph (Kwak, *et al.*) with 61.6 million vertices and 1.47 billion edges in 105 minutes. We use GraphCT to analyze public data from Twitter, a microblogging network. Twitter's message connections appear primarily tree-structured as a news dissemination system. Within the

involves over 400 million active users with an ave 120 'friendship' connections each and sharing 5 references to items each month [11].

One analysis approach treats the interactions as and applies tools from graph theory, social r analysis, and scale-free networks [29]. Howev volume of data that must be processed to appl techniques overwhelms current computational capa Even well-understood analytic methodologies advances in both hardware and software to proc growing corpus of social media.

Social media provides staggering amounts

### Top 15 Users by Betweenness Centrality

| Rank | Data Set | |
| --- | --- | --- |
| | H1N1 | atlflood |
| 1 | @CDCFlu | @ajc |
| 2 | @addthis | @driveafaste |
| 3 | @Official_PAX | @ATLCheap |
| 4 | @FluGov | @TWCi |
| 5 | @nytimes | @HelloNorthGA |
| 6 | @tweetmeme | @11AliveNews |
| 7 | @mercola | @WSB_TV |
| 8 | @CNN | @shaunking |
| 9 | @backstreetboys | @Carl |
| 10 | @EllieSmith_x | @SpaceyG |
| 11 | @TIME | @ATLINtownPa |
| 12 | @CDCemergency | @TJsDJs |
| 13 | @CDC_eHealth | @ATLien |
| 14 | @perezhilton | @MarshallRamsey |
| 15 | @billmaher | @Kanye |

H1N1
17k vertices

1184 vertices

Fig. 3  Subcommunity filtering on Twitter data sets

Image credit: bioethicsinstitute.org

David A. Bader

# STING Initiative: Focusing on Globally Significant Grand Challenges

- Many globally-significant grand challenges can be modeled by **Spatio-Temporal Interaction Networks and Graphs** (or "STING").

- Emerging real-world graph problems include:
  - Detecting community structure in large social networks
  - Defending the nation against cyber-based attacks
  - Discovering insider threats (e.g. Ft. Hood shooter, WikiLeaks)
  - Improving the resilience of the electric power grid
  - Detecting and preventing disease in human populations.

- Unlike traditional applications in computational science and engineering, solving these problems at scale often raises new research challenges due to:
  - Sparsity and the lack of locality in the massive data
  - Design of parallel algorithms for massive, streaming data analytics
  - The need for new exascale supercomputers that are energy-efficient, resilient, and easy-to-program

# STINGER – Time Frame



**Pre-1999**

Streaming graph need arises (over a decade ago)

STINGER is officially proposed. May 2009

**2009**

First prototype, clustering coefficients. Apr 2010

**2010**

Structure tracking of streaming social networks. Apr 2011

**2011**

Dynamic betweenness centrality algorithm. Sep 2012

High Performance Data Structure for Streaming Graphs. Sep 2012
**HPEC BEST PAPER AWARD**

**2012**

Streaming connected component, Dec 2013

**2013**

Performance evaluation of open-source graph data-bases. Feb 2014

**2014**

Community detection in dynamic networks. Sep 2015

**2015**

PageRank for Streaming Graphs. May 2016

**2016**

# STING Extensible Representation (STINGER)

## *Design goals*

- Enable algorithm designers to implement dynamic graph algorithms with ease.

- Portable semantics for various platforms

- Good performance for all types of graph problems and algorithms - static and dynamic.

- Assumes globally addressable memory access

- Support multiple, parallel readers and a single writer
  - One server manages the graph data structures
  - Multiple analytics run in background with read-only permissions.

# STING Extensible Representation (STINGER)

- Semi-dense edge list blocks with free space

- Compactly stores timestamps, types, weights

- Maps from application IDs to storage IDs

- Deletion by negating IDs, separate compaction

# STINGER as an analysis package

**Anything that a static graph package can do (and a whole lot more):**

**Parallel agglomerative clustering:**
Find clusters that are optimized for a user-defined edge scoring function.

**K-core Extraction:**
Extract additional communities and filter noisy high-degree vertices.

**Classic breadth-first search:**
Performs a parallel breadth-first search of the graph starting at a given source vertex to find shortest paths.

**Parallel connected components:**
Finds the connected components in a static network.

**AND…**

**Streaming edge insertions and deletions:**
New edge insertions, updates, and deletions in batches or individually.
Optimized to update at rates of over 3 million edges per second on graphs of one billion edges.

**Streaming clustering coefficients:**
Tracks the local and global clustering coefficients of a graph.

**Streaming connected components:**
Real time tracking of the connected components.

**Streaming Betweenness Centrality:**
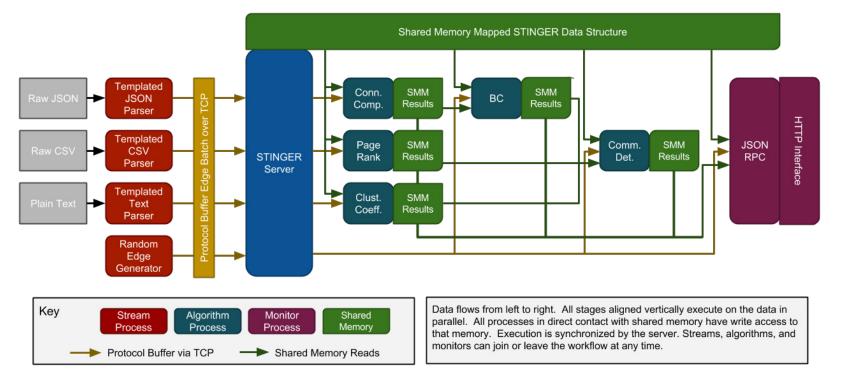Find the key points within information flows and structural vulnerabilities.

**Streaming community detection:**
Track and update the community structures within the graph as they change.

# STING: High-level architecture



Shared Memory Mapped STINGER Data Structure

Raw JSON → Templated JSON Parser
Raw CSV → Templated CSV Parser
Plain Text → Templated Text Parser
Random Edge Generator

Protocol Buffer Edge Batch over TCP → STINGER Server

Conn. Comp. | SMM Results
Page Rank | SMM Results
Clust. Coeff. | SMM Results

BC | SMM Results
Comm. Det. | SMM Results

JSON RPC | HTTP Interface

**Key**

| Stream Process | Algorithm Process | Monitor Process | Shared Memory |

→ Protocol Buffer via TCP    → Shared Memory Reads

Data flows from left to right. All stages aligned vertically execute on the data in parallel. All processes in direct contact with shared memory have write access to that memory. Execution is synchronized by the server. Streams, algorithms, and monitors can join or leave the workflow at any time.

⚠ Server: Graph storage, kernel orchestration
⚠ OpenMP + sufficiently POSIX-ish
⚠ Multiple processes for resilience

10 September 2019

David A. Bader

# STINGER Summary

- Massive-Scale Streaming Analytics require
  - Simple programming model
    - Simple API.
    - CSR-like in concept.
    - STINGER has a lot more under the hood.
  - Extremely fast updates
    - Millions of updates per second.
    - These must not be bottlenecks for updating an analytic.
    - STINGER offers these

- STINGER has major performance benefits
  - Thousands of times faster than static graph computation.
  - Hundreds of thousands of updates per second for numerous analytics.
  - Real-time monitoring of underlying network.

David A. Bader

# Graphs for Analysts: Putting it all together

- Today homeland security knowledge is captured in *reports*
  - These are written in a natural language, can be multiple pages, include references, footnotes, figures, etc.

- Analysts manually discover and re-discover relationships among the unclassified security reports and open data
  - Searching and filtering help uncover relationships, but all complex queries must be done by hand and manually

- Graphs are a natural and important structure for connecting this knowledge

# Approach

- We developed a *knowledge graph* for the Department of Homeland Security, with a focus on *special events*, working with DHS/FBI/NYPD

- This knowledge graph is built on top of:
  - Reports
  - Spreadsheets
  - Raw datasets

- Enterprise-class performance with STINGER

NJIT
New Jersey Institute
of Technology

# Using the Knowledge Graph

After the graph has been created, **relationships** between reports and events can be maintained

Analysts are able to use this graph to quickly triage new events / information and **link together** otherwise hard-to-connect data

Furthermore, **advanced and predictive analytics** can be run to gain further insight

**New Jersey Institute of Technology**

# Data Sources

Several sources contribute to homeland security knowledge graph

Focus on three datasets:

- A collection of reports provided by partners
- The Global Terrorism Database (GTD) from the University of Maryland
- The Homegrown Violent Extremists (HVE) list from the University of Nebraska Omaha

NJIT
New Jersey Institute
of Technology

# Ingesting Unstructured Text

**Vertices**: Proper nouns were extracted and given labels of *Media, Events, Groups, Individuals,* and *Locations*

**Properties**: Attacks and threats were extracted and placed into the *AttackType* label (e.g. explosive devices or vehicle attacks)
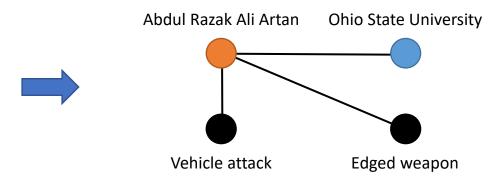
**Edges** were created when a report indicated a connection between the vertices

NJIT
New Jersey Institute of Technology

# Parsing Example

On November 28, 2016, a terrorist **vehicle-ramming** and **stabbing attack** occurred at 9:52 a.m. EST at **Ohio State University** (OSU)'s Watts Hall in Columbus, Ohio. The attacker, Somali refugee **Abdul Razak Ali Artan**, was shot and killed by the first responding OSU police officer, and 13 people were hospitalized for injuries.
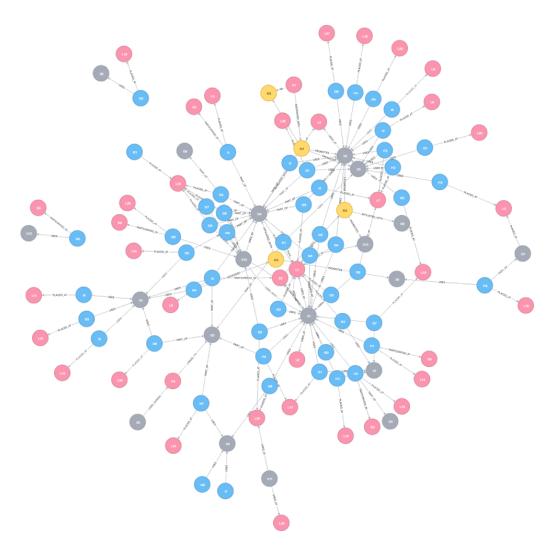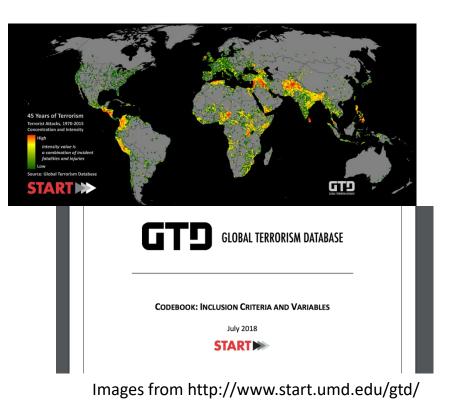
*(source: Wikipedia)*

# Ingested Data



- 29 Joint Threat Assessments reports were parsed
    - Provided by DHS
    - Recent (2016 to 2018) and relevant to the problem

- 114 vertices, 163 edges

NJIT
New Jersey Institute of Technology
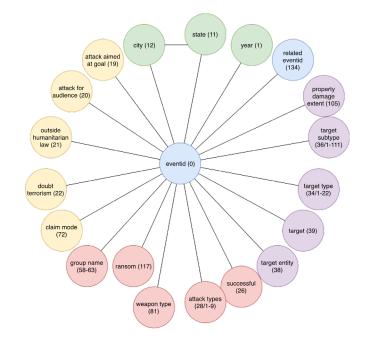
# Ingesting the Global Terrorism Database / HVE

- Database is provided in the form of a spreadsheets with structured coding

- Converted the spreadsheets to CSV files, and then ingested primarily using `LOAD CSV`

- Pre-processing performed with Python scripts



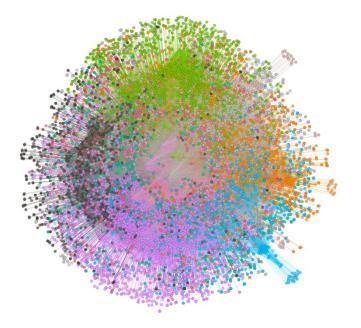Images from http://www.start.umd.edu/gtd/

# Graph Schema and Ingestion

- Schema is focused around database rows
  - This creates a nearly bipartite graph from rows and columns
- Other connections are then created through queries
- 325K vertices, 2.7M edges

# Running Analytics on the Homeland Security Knowledge Graph



Louvain running on US events, visualized in Gephi

Node Importance

| Id | eigencentrality | pageranks |
|---|---|---|
| 4 | 1 | 0.024539 |
| 6 | 0.998713 | 0.024502 |
| 5 | 0.99301 | 0.024298 |
| 8 | 0.846495 | 0.020458 |
| 7 | 0.829617 | 0.019987 |
| 28 | 0.529627 | 0.012339 |
| 22 | 0.525788 | 0.012185 |
| 16 | 0.479301 | 0.011143 |
| 32 | 0.306425 | 0.007603 |
| 39 | 0.301179 | 0.007432 |
| 67 | 0.289466 | 0.006846 |
| 18 | 0.219309 | 0.005264 |
| 27 | 0.198742 | 0.005009 |

| Eigencentrality | Group |
|---|---|
| 0.198742 | Unknown |
| 0.073040 | Anti-Abortion extremists |
| 0.063670 | Left-Wing Militants |
| 0.047164 | FALN |
| 0.032517 | New World Liberation Front |
| 0.030430 | Black Nationalists |
| 0.028223 | White extremists |
| 0.027993 | Animal Liberation Front |
| … | … |

NJIT
New Jersey Institute
of Technology

# Conclusions

- Solving massive-scale analytics will require new
  - High-performance computing platforms
  - Streaming algorithms
  - Energy-efficient implementations
- Mapping applications to high performance architectures may yield performance improvements of six or more orders of magnitude.
- Solving real-world challenges such as:
  - Urban sustainability
  - Healthcare analytics
  - Trustworthy, Free and Fair Elections
  - Insider threat detection
  - Utility infrastructure protection
  - Cyberattack defense
  - Disease outbreak and epidemic monitoring

NJIT
New Jersey Institute
of Technology

# Acknowledgments

- Jason Riedy, Research Scientist, (Georgia Tech)

- Oded Green, Research Scientist, (Georgia Tech)

- Current Graduate Students (Georgia Tech):
  - **Xiaojing An**
  - **James Fox**
  - **Kasimir Gabert**
  - **Euna Kim**

- Recent Bader Alumni:
  - **Dr. Eisha Nathan** (Lawrence Livermore National Lab)
  - **Dr. Vipin Sachdeva** (IBM)
  - **Dr. Anita Zakrzewska** (Lawrence Livermore National Lab)
  - **Dr. Lluis Miquel Munguia** (Google)
  - **Prof. Kamesh Madduri** (Penn State)
  - **Dr. David Ediger** (GTRI)
  - **Dr. James Fairbanks** (GTRI)
  - **Dr. Seunghwa Kang** (Pacific Northwest National Lab)

# Contribution

- Developed new algorithm for local community detection that incrementally updates results

- Outputs results similar to static re-computation
  - Average recall and precision 0.80-0.99 and 0.59-0.98 across real graphs
  - Recall and precision do not decrease over time
  - Can start with no data and build community

- Faster than re-computing with static algorithm
  - Dynamic algorithm most beneficial for small batch sizes
  - Reaches up to two orders of magnitude dynamic speedup

# Bader, Related Recent Publications (2005-2009)

- D.A. Bader, G. Cong, and J. Feo, "**On the Architectural Requirements for Efficient Execution of Graph Algorithms**," *The 34th International Conference on Parallel Processing* (ICPP 2005), pp. 547-556, Georg Sverdrups House, University of Oslo, Norway, June 14-17, 2005.

- D.A. Bader and K. Madduri, "**Design and Implementation of the HPCS Graph Analysis Benchmark on Symmetric Multiprocessors**," *The 12th International Conference on High Performance Computing* (HiPC 2005), D.A. Bader *et al.*, (eds.), Springer-Verlag LNCS 3769, 465-476, Goa, India, December 2005.

- D.A. Bader and K. Madduri, "**Designing Multithreaded Algorithms for Breadth-First Search and st-connectivity on the Cray MTA-2**," *The 35th International Conference on Parallel Processing* (ICPP 2006), Columbus, OH, August 14-18, 2006.

- D.A. Bader and K. Madduri, "**Parallel Algorithms for Evaluating Centrality Indices in Real-world Networks**," *The 35th International Conference on Parallel Processing* (ICPP 2006), Columbus, OH, August 14-18, 2006.

- K. Madduri, D.A. Bader, J.W. Berry, and J.R. Crobak, "**Parallel Shortest Path Algorithms for Solving Large-Scale Instances**," *9th DIMACS Implementation Challenge -- The Shortest Path Problem*, DIMACS Center, Rutgers University, Piscataway, NJ, November 13-14, 2006.

- K. Madduri, D.A. Bader, J.W. Berry, and J.R. Crobak, "**An Experimental Study of A Parallel Shortest Path Algorithm for Solving Large-Scale Graph Instances**," *Workshop on Algorithm Engineering and Experiments* (ALENEX), New Orleans, LA, January 6, 2007.

- J.R. Crobak, J.W. Berry, K. Madduri, and D.A. Bader, "**Advanced Shortest Path Algorithms on a Massively-Multithreaded Architecture**," *First Workshop on Multithreaded Architectures and Applications* (MTAAP), Long Beach, CA, March 30, 2007.

- D.A. Bader and K. Madduri, "**High-Performance Combinatorial Techniques for Analyzing Massive Dynamic Interaction Networks**," *DIMACS Workshop on Computational Methods for Dynamic Interaction Networks*, DIMACS Center, Rutgers University, Piscataway, NJ, September 24-25, 2007.

- D.A. Bader, S. Kintali, K. Madduri, and M. Mihail, "**Approximating Betewenness Centrality**," The *5th Workshop on Algorithms and Models for the Web-Graph* (WAW2007), San Diego, CA, December 11-12, 2007.

- David A. Bader, Kamesh Madduri, Guojing Cong, and John Feo, "**Design of Multithreaded Algorithms for Combinatorial Problems**," in S. Rajasekaran and J. Reif, editors, *Handbook of Parallel Computing: Models, Algorithms, and Applications*, CRC Press, Chapter 31, 2007.

- Kamesh Madduri, David A. Bader, Jonathan W. Berry, Joseph R. Crobak, and Bruce A. Hendrickson, "**Multithreaded Algorithms for Processing Massive Graphs**," in D.A. Bader, editor, *Petascale Computing: Algorithms and Applications*, Chapman & Hall / CRC Press, Chapter 12, 2007.

- D.A. Bader and K. Madduri, "**SNAP, Small-world Network Analysis and Partitioning: an open-source parallel graph framework for the exploration of large-scale networks**," *22nd IEEE International Parallel and Distributed Processing Symposium* (IPDPS), Miami, FL, April 14-18, 2008.

- S. Kang, D.A. Bader, "**An Efficient Transactional Memory Algorithm for Computing Minimum Spanning Forest of Sparse Graphs**," 14th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP), Raleigh, NC, February 2009.

- Karl Jiang, David Ediger, and David A. Bader. "**Generalizing k-Betweenness Centrality Using Short Paths and a Parallel Multithreaded Implementation**." The 38th International Conference on Parallel Processing (ICPP), Vienna, Austria, September 2009.

- Kamesh Madduri, David Ediger, Karl Jiang, David A. Bader, Daniel Chavarría-Miranda. "**A Faster Parallel Algorithm and Efficient Multithreaded Implementations for Evaluating Betweenness Centrality on Massive Datasets**." 3rd Workshop on Multithreaded Architectures and Applications (MTAAP), Rome, Italy, May 2009.

- David A. Bader, et al. "**STINGER: Spatio-Temporal Interaction Networks and Graphs (STING) Extensible Representation**." 2009.

# Bader, Related Recent Publications (2010-2011)

- David Ediger, Karl Jiang, E. Jason Riedy, and David A. Bader. "**Massive Streaming Data Analytics: A Case Study with Clustering Coefficients,**" Fourth Workshop in Multithreaded Architectures and Applications (MTAAP), Atlanta, GA, April 2010.

- Seunghwa Kang, David A. Bader. "**Large Scale Complex Network Analysis using the Hybrid Combination of a MapReduce cluster and a Highly Multithreaded System:**," Fourth Workshop in Multithreaded Architectures and Applications (MTAAP), Atlanta, GA, April 2010.

- David Ediger, Karl Jiang, Jason Riedy, David A. Bader, Courtney Corley, Rob Farber and William N. Reynolds. "**Massive Social Network Analysis: Mining Twitter for Social Good,**" The 39th International Conference on Parallel Processing (ICPP 2010), San Diego, CA, September 2010.

- Virat Agarwal, Fabrizio Petrini, Davide Pasetto and David A. Bader. "**Scalable Graph Exploration on Multicore Processors,**" *The 22nd IEEE and ACM Supercomputing Conference* (SC10), New Orleans, LA, November 2010.

- Z. Du, Z. Yin, W. Liu, and D.A. Bader, "**On Accelerating Iterative Algorithms with CUDA: A Case Study on Conditional Random Fields Training Algorithm for Biological Sequence Alignment,**" IEEE International Conference on Bioinformatics & Biomedicine, Workshop on Data-Mining of Next Generation Sequencing Data (NGS2010), Hong Kong, December 20, 2010.

- D. Ediger, J. Riedy, H. Meyerhenke, and D.A. Bader, "**Tracking Structure of Streaming Social Networks,**" 5th Workshop on Multithreaded Architectures and Applications (MTAAP), Anchorage, AK, May 20, 2011.

- D. Mizell, D.A. Bader, E.L. Goodman, and D.J. Haglin, "**Semantic Databases and Supercomputers,**" 2011 Semantic Technology Conference (SemTech), San Francisco, CA, June 5-9, 2011.

- P. Pande and D.A. Bader, "**Computing Betweenness Centrality for Small World Networks on a GPU,**" *The 15th Annual High Performance Embedded Computing Workshop* (HPEC), Lexington, MA, September 21-22, 2011.

- David A. Bader, Christine Heitsch, and Kamesh Madduri, "**Large-Scale Network Analysis,**" in J. Kepner and J. Gilbert, editor, *Graph Algorithms in the Language of Linear Algebra*, SIAM Press, Chapter 12, pages 253-285, 2011.

- Jeremy Kepner, David A. Bader, Robert Bond, Nadya Bliss, Christos Faloutsos, Bruce Hendrickson, John Gilbert, and Eric Robinson, "**Fundamental Questions in the Analysis of Large Graphs,**" in J. Kepner and J. Gilbert, editor, *Graph Algorithms in the Language of Linear Algebra*, SIAM Press, Chapter 16, pages 353-357, 2011.

# Bader, Related Recent Publications (2012)

- E.J. Riedy, H. Meyerhenke, D. Ediger, and D.A. Bader, "**Parallel Community Detection for Massive Graphs**," The 9th International Conference on Parallel Processing and Applied Mathematics (PPAM 2011), Torun, Poland, September 11-14, 2011. Lecture Notes in Computer Science, 7203:286-296, 2012.

- E.J. Riedy, D. Ediger, D.A. Bader, and H. Meyerhenke, "**Parallel Community Detection for Massive Graphs**," 10th DIMACS Implementation Challenge -- Graph Partitioning and Graph Clustering, Atlanta, GA, February 13-14, 2012.

- E.J. Riedy, H. Meyerhenke, D.A. Bader, D. Ediger, and T. Mattson, "**Analysis of Streaming Social Networks and Graphs on Multicore Architectures**," The 37th IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP), Kyoto, Japan, March 25-30, 2012.

- J. Riedy, H. Meyerhenke, and D.A. Bader, "**Scalable Multi-threaded Community Detection in Social Networks**," 6th Workshop on Multithreaded Architectures and Applications (MTAAP), Shanghai, China, May 25, 2012.

- H. Meyerhenke, E.J. Riedy, and D.A. Bader, "**Parallel Community Detection in Streaming Graphs**," Minisymposium on Parallel Analysis of Massive Social Networks, *15th SIAM Conference on Parallel Processing for Scientific Computing* (PP12), Savannah, GA, February 15-17, 2012.

- D. Ediger, E.J. Riedy, H. Meyerhenke, and D.A. Bader, "**Analyzing Massive Networks with GraphCT**," Poster Session, *15th SIAM Conference on Parallel Processing for Scientific Computing* (PP12), Savannah, GA, February 15-17, 2012.

- R.C. McColl, D. Ediger, and D.A. Bader, "**Many-Core Memory Hierarchies and Parallel Graph Analysis**," Poster Session, *15th SIAM Conference on Parallel Processing for Scientific Computing* (PP12), Savannah, GA, February 15-17, 2012.

- E.J. Riedy, D. Ediger, H. Meyerhenke, and D.A. Bader, "**STING: Software for Analysis of Spatio-Temporal Interaction Networks and Graphs,**" Poster Session, *15th SIAM Conference on Parallel Processing for Scientific Computing* (PP12), Savannah, GA, February 15-17, 2012.

- Y. Chai, Z. Du, D.A. Bader, and X. Qin, "**Efficient Data Migration to Conserve Energy in Streaming Media Storage Systems,**" *IEEE Transactions on Parallel & Distributed Systems*, 2012.

- M. S. Swenson, J. Anderson, A. Ash, P. Gaurav, Z. Sükösd, D.A. Bader, S.C. Harvey and C.E Heitsch, "**GTfold: Enabling parallel RNA secondary structure prediction on multi-core desktops,**" *BMC Research Notes*, 5:341, 2012.

- D. Ediger, K. Jiang, E.J. Riedy, and D.A. Bader, "**GraphCT: Multithreaded Algorithms for Massive Graph Analysis**," *IEEE Transactions on Parallel & Distributed Systems*, 2012.

- D.A. Bader and K. Madduri, "**Computational Challenges in Emerging Combinatorial Scientific Computing Applications**," in O. Schenk, editor, *Combinatorial Scientific Computing*, Chapman & Hall / CRC Press, Chapter 17, pages 471-494, 2012.

- O. Green, R. McColl, and D.A. Bader, "**GPU Merge Path -- A GPU Merging Algorithm**," *26th ACM International Conference on Supercomputing* (ICS), San Servolo Island, Venice, Italy, June 25-29, 2012.

- O. Green, R. McColl, and D.A. Bader, "**A Fast Algorithm for Streaming Betweenness Centrality**," *4th ASE/IEEE International Conference on Social Computing* (SocialCom), Amsterdam, The Netherlands, September 3-5, 2012.

- D. Ediger, R. McColl, J. Riedy, and D.A. Bader, "**STINGER: High Performance Data Structure for Streaming Graphs**," *The IEEE High Performance Extreme Computing Conference* (HPEC), Waltham, MA, September 20-22, 2012. **Best Paper Award.**

- J. Marandola, S. Louise, L. Cudennec, J.-T. Acquaviva and D.A. Bader, "**Enhancing Cache Coherent Architecture with Access Patterns for Embedded Manycore Systems**," *14th IEEE International Symposium on System-on-Chip* (SoC), Tampere, Finland, October 11-12, 2012.

- L.M. Munguía, E. Ayguade, and D.A. Bader, "**Task-based Parallel Breadth-First Search in Heterogeneous Environments**," *The 19th Annual IEEE International Conference on High Performance Computing* (HiPC), Pune, India, December 18-21, 2012.

# Bader, Related Recent Publications (2013)

- X. Liu, P. Pande, H. Meyerhenke, and D.A. Bader, "**PASQUAL: Parallel Techniques for Next Generation Genome Sequence Assembly,**" *IEEE Transactions on Parallel & Distributed Systems*, 24(5):977-986, 2013.

- David A. Bader, Henning Meyerhenke, Peter Sanders, and Dorothea Wagner (eds.), ***Graph Partitioning and Graph Clustering***, American Mathematical Society, 2013.

- E. Jason Riedy, Henning Meyerhenke, David Ediger and David A. Bader, "**Parallel Community Detection for Massive Graphs**," in David A. Bader, Henning Meyerhenke, Peter Sanders, and Dorothea Wagner (eds.), *Graph Partitioning and Graph Clustering*, American Mathematical Society, Chapter 14, pages 207-222, 2013.

- S. Kang, D.A. Bader, and R. Vuduc, "**Energy-Efficient Scheduling for Best-Effort Interactive Services to Achieve High Response Quality**," *27th IEEE International Parallel and Distributed Processing Symposium* (IPDPS), Boston, MA, May 20-24, 2013.

- J. Riedy and D.A. Bader, "**Multithreaded Community Monitoring for Massive Streaming Graph Data**," *7th Workshop on Multithreaded Architectures and Applications* (MTAAP), Boston, MA, May 24, 2013.

- D. Ediger and D.A. Bader, "**Investigating Graph Algorithms in the BSP Model on the Cray XMT**," *7th Workshop on Multithreaded Architectures and Applications* (MTAAP), Boston, MA, May 24, 2013.

- O. Green and D.A. Bader, "**Faster Betweenness Centrality Based on Data Structure Experimentation**," *International Conference on Computational Science* (ICCS), Barcelona, Spain, June 5-7, 2013.

- Z. Yin, J. Tang, S. Schaeffer, and D.A. Bader, "**Streaming Breakpoint Graph Analytics for Accelerating and Parallelizing the Computation of DCJ Median of Three Genomes**," *International Conference on Computational Science* (ICCS), Barcelona, Spain, June 5-7, 2013.

- T. Senator, D.A. Bader, et al., "**Detecting Insider Threats in a Real Corporate Database of Computer Usage Activities**," *19th ACM SIGKDD Conference on Knowledge Discovery and Data Mining* (KDD), Chicago, IL, August 11-14, 2013.

- J. Fairbanks, D. Ediger, R. McColl, D.A. Bader and E. Gilbert, "**A Statistical Framework for Streaming Graph Analysis**," *IEEE/ACM International Conference on Advances in Social Networks Analysis and Modeling* (ASONAM), Niagara Falls, Canada, August 25-28, 2013.

- A. Zakrzewska and D.A. Bader, "**Measuring the Sensitivity of Graph Metrics to Missing Data**," *10th International Conference on Parallel Processing and Applied Mathematics* (PPAM), Warsaw, Poland, September 8-11, 2013.

- O. Green and D.A. Bader, "**A Fast Algorithm for Streaming Betweenness Centrality**," *5th ASE/IEEE International Conference on Social Computing* (SocialCom), Washington, DC, September 8-14, 2013.

- R. McColl, O. Green, and D.A. Bader, "**A New Parallel Algorithm for Connected Components in Dynamic Graphs**," *The 20th Annual IEEE International Conference on High Performance Computing* (HiPC), Bangalore, India, December 18-21, 2013.

# Bader, Related Recent Publications (2014-2015)

- R. McColl, D. Ediger, J. Poovey, D. Campbell, and D.A. Bader, "**A Performance Evaluation of Open Source Graph Databases**," *The 1st Workshop on Parallel Programming for Analytics Applications* (PPAA 2014) held in conjunction with the *19th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming* (PPoPP 2014), Orlando, Florida, February 16, 2014.

- O. Green, L.M. Munguia, and D.A. Bader, "**Load Balanced Clustering Coefficients**," *The 1st Workshop on Parallel Programming for Analytics Applications* (PPAA 2014) held in conjunction with the *19th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming* (PPoPP 2014), Orlando, Florida, February 16, 2014.

- A. McLaughlin and D.A. Bader, "**Revisiting Edge and Node Parallelism for Dynamic GPU Graph Analytics**," *8th Workshop on Multithreaded Architectures and Applications* (MTAAP), held in conjunction with *The IEEE International Parallel and Distributed Processing Symposium (IPDPS 2014)*, Phoenix, AZ, May 23, 2014.

- Z. Yin, J. Tang, S. Schaeffer, D.A. Bader, "**A Lin-Kernighan Heuristic for the DCJ Median Problem of Genomes with Unequal Contents**," *20th International Computing and Combinatorics Conference* (COCOON), Atlanta, GA, August 4-6, 2014.

- Y. You, D.A. Bader and M.M. Dehnavi, "**Designing an Adaptive Cross-Architecture Combination for Graph Traversal,**" *The 43rd International Conference on Parallel Processing* (ICPP 2014), Minneapolis, MN, September 9-12, 2014.

- A. McLaughlin, J. Riedy, and D.A. Bader, "**Optimizing Energy Consumption and Parallel Performance for Betweenness Centrality using GPUs**," *The 18th Annual IEEE High Performance Extreme Computing Conference* (HPEC), Waltham, MA, September 9-11, 2014.

- A. McLaughlin and D.A. Bader, "**Scalable and High Performance Betweenness Centrality on the GPU**," *The 26th IEEE and ACM Supercomputing Conference* (SC14), New Orleans, LA, November 16-21, 2014. **Best Student Paper Finalist.**

- D. Dauwe, E. Jonardi, R. Friese, S. Pasricha, A.A. Maciejewski, D.A. Bader, and H.J. Siegel, "**A Methodology for Co-Location Aware Application Performance Modeling in Multicore Computing,**" 17th Workshop on Advances on Parallel and Distributed Processing Symposium (APDCM), Hyderabad, India, May 25, 2015.

- A. Zakrzewska and D.A. Bader, "**Fast Incremental Community Detection on Dynamic Graphs**," 11th International Conference on Parallel Processing and Applied Mathematics (PPAM), Krakow, Poland, September 6-9, 2015.

- A. McLaughlin, J. Riedy, and D.A. Bader, "**An Energy-Efficient Abstraction for Simultaneous Breadth-First Searches**," The 19th Annual IEEE High Performance Extreme Computing Conference (HPEC), Waltham, MA, September 15-17, 2015.

- A. McLaughlin, D. Merrill, M. Garland and D.A. Bader, "**Parallel Methods for Verifying the Consistency of Weakly-Ordered Architectures**," The 24th International Conference on Parallel Architectures and Compilation Techniques (PACT), San Francisco, CA, October 18-21, 2015.

- A. McLaughlin and D.A. Bader, "**Fast Execution of Simultaneous Breadth-First Searches on Sparse Graphs**,'' The 21st IEEE International Conference on Parallel and Distributed Systems (ICPADS), Melbourne, Australia, December 14-17, 2015.

David A. Bader

# Bader, Related Recent Publications (2016-2017)

- David Bader, Aleksandra Michalewicz, Oded Green, Jessie Birkett-Rees, Jason Riedy, James Fairbanks, and Anita Zakrzewska, "*Semantic database applications at the Samtavro Cemetery, Georgia*," , The 44th Computer Applications and Quantitative Methods in Archaeology Conference (CAA), Oslo, Norway, March 29 – April 2, 2016.

- Vipin Sachdeva, Srinivas Aluru, David A. Bader, "*A Memory and Time Scalable Parallelization of the Reptile Error-Correction Code*," 15th IEEE International Workshop on High Performance Computational Biology (HiCOMB), Chicago, IL, May 23, 2016.

- James Fairbanks, Anita Zakrzewska, and David A. Bader, "*New Stopping Criteria For Spectral Partitioning*," IEEE/ACM International Conference on Advances in Social Networks Analysis and Modeling (ASONAM), San Francisco, CA, August 18-21, 2016.

- Anita Zakrzewska, Eisha Nathan, James Fairbanks, and David A. Bader, "*A Local Measure of Community Change in Dynamic Graphs*," IEEE/ACM International Conference on Advances in Social Networks Analysis and Modeling (ASONAM), San Francisco, CA, August 18-21, 2016.

- Anita Zakrzewska and David A. Bader, "*Aging Data in Dynamic Graphs: A Comparative Study*," 2nd International Workshop on Dynamics in Networks (DyNo), held in conjunction with IEEE/ACM International Conference on Advances in Social Networks Analysis and Modeling (ASONAM), San Francisco, CA, August 18, 2016.

- O. Green and D.A. Bader, "*cuSTINGER: Supporting Dynamic Graph Algorithms for GPUs,*" The 20th Annual IEEE High Performance Extreme Computing Conference (HPEC), Waltham, MA, September 13-15, 2016.

- Jeremy Kepner, Peter Aaltonen, David A. Bader, Aydin Buluc, Franz Franchetti, John Gilbert, Dylan Hutchison, Manoj Kumar, Andrew Lumsdaine, Henning Meyerhenke, Scott McMillan, Jose Moreira, John D. Owens, Carl Yang, Marcin Zalewski, and Timothy Mattson, "*Mathematical Foundations of the GraphBLAS*," The 20th Annual IEEE High Performance Extreme Computing Conference (HPEC), Waltham, MA, September 13-15, 2016.

- X. Hui, Z. Du, J. Liu, H. Sun, Y. He and D.A. Bader, "*When Good Enough Is Better: Energy-Aware Scheduling for Multicore Servers,*" 13th Workshop on High-Performance, PowerAware Computing (HPPAC), Orlando, FL, May 29, 2017.

- E. Nathan, G. Sanders, J. Fairbanks, V. Henson and D.A. Bader, "*Graph Ranking Guarantees for Numerical Approximations to Katz Centrality*," International Conference on Computational Science (ICCS), Zurich, Switzerland, June 12-14, 2017.

- Anita Zakrzewska and David A. Bader, "*Streaming Graph Sampling with Size Restrictions*," IEEE/ACM International Conference on Advances in Social Networks Analysis and Modeling (ASONAM), Sydney, Australia, July 31 - August 3, 2017.

- Eisha Nathan and David A. Bader, "*A Dynamic Algorithm for Updating Katz Centrality in Graphs*," , IEEE/ACM International Conference on Advances in Social Networks Analysis and Modeling (ASONAM), Sydney, Australia, July 31 - August 3, 2017.

- E. Nathan and D.A. Bader, "*Approximating Personalized Katz Centrality in Dynamic Graphs*," , 12th International Conference on Parallel Processing and Applied Mathematics(PPAM), Lublin, Poland, September 10-13, 2017.

- O. Green, J. Fox, E. Kim, F. Busato, N. Bombieri, K. Lakhotia, S. Zhou, S. Singapura, H. Zeng, R. Kannan, V. Prasanna, D. Bader, "*Quickly Finding a Truss in a Haystack*", IEEE High Performance Extreme Computing Conference (HPEC), Waltham, Massachusetts, 2017 **(HPEC Graph Challenge Innovation Award)**

- S. Zhou, K. Lakhotia, S. Singapura, H. Zeng, R. Kannan, V. Prasanna, J. Fox, E. Kim, O. Green, D. Bader, "*Design and Implementation of Parallel PageRank on Multicore Platforms*", IEEE High Performance Extreme Computing Conference (HPEC), Waltham, Massachusetts, 2017 **(HPEC Graph Challenge Student Innovation Award)**

- D. Makkar, D. Bader, O. Green, "*Deterministic and Parallel Triangle Counting in Streaming Graphs*", IEEE International Conference on High Performance Computing, Data, and Analytics, Jaipur, India, 2017

# Bader, Related Recent Publications (2018)

- *E. Nathan, G. Sanders, V. Emden Henson, D.A. Bader, ``**Numerically Approximating Centrality for Graph Ranking Guarantees**,'' Journal of Computational Science, 26:205-216, May 2018.*
- *A. McLaughlin and D.A. Bader, ``**Accelerating GPU Betweenness Centrality**,'' Communications of the ACM, 61(8):85-92, August 2018.*
- *Eisha Nathan, James Fairbanks, and David A. Bader, ``**Ranking in Dynamic Graphs using Exponential Centrality**,'' 6th International Conference on Complex Networks and Their Applications, Lyon, France, November 29 - December 1, 2017. In: Cherifi C., Cherifi H., Karsai M., Musolesi M. (eds), Complex Networks & Their Applications VI, Complex Networks 2016 2017. Springer Studies in Computational Intelligence, 689:378-389, 2018.*
- *Devavret Makkar, David A. Bader, and Oded Green, ``**Exact and Parallel Triangle Counting in Dynamic Graphs**,'' The 24th IEEE International Conference on High Performance Computing (HiPC), Jaipur, India, December 18-21, 2017.*
- *D.A. Bader, ``**Massive-Scale Streaming Analytics: Models, Parallelsim, and Real-World Applications**,'' The 30th ACM Symposium on Parallelism in Algorithms and Architectures (SPAA), Vienna, Austria, p. 193, July 16-18, 2018.*
- *Alexander van der Grinten, Elisabetta Bergamini, Oded Green, David A. Bader, and Henning Meyerhenke, ``**Scalable Katz Ranking Computation in Large Dynamic Graphs**,'' The 26th Annual European Symposium on Algorithms (ESA), Helsinki, Finland, August 20-22, 2018.*
- *Federico Busato, Oded Green, Nicola Bombieri, and David A. Bader, ``**Hornet: An Efficient Data Structure for Dynamic Sparse Graphs and Matrices on GPUs,**'' The 22nd Annual IEEE High Performance Extreme Computing Conference (HPEC), Waltham, MA, September 25-27, 2018.*
- *James Fox, Oded Green, Kasimir Gabert, Xiaojing An, and David A. Bader, ``**Fast and Adaptive List Intersections on the GPU**,'' The 22nd Annual IEEE High Performance Extreme Computing Conference (HPEC), Waltham, MA, September 25-27, 2018. **Graph Challenge Finalist.***
- *Oded Green, James Fox, Alex Watkins, Alok Tripathy, Kasimir Gabert, Euna Kim, An Xiaojing, Kumar Aatish, and David A. Bader, ``**Logarithmic Radix Binning and Vectorized Triangle Counting**,'' The 22nd Annual IEEE High Performance Extreme Computing Conference (HPEC), Waltham, MA, September 25-27, 2018. **Graph Challenge Innovation Award.***
- *D.A. Bader, ``**Massive-Scale Streaming Analytics**,'' Minisymposium on Scalable and Dynamic Graph Algorithms, 17th SIAM Conference on Parallel Processing for Scientific Computing (PP18), Tokyo, Japan, March 7-10, 2018.*
- *D.A. Bader, ``**Tracking Communities in Streaming Graphs**,'' Minisymposium on High Performance Graph Algorithms, The 5th ACM Platform for Advanced Scientific Computing Conference (PASC), Basel, Switzerland, July 3, 2018.*