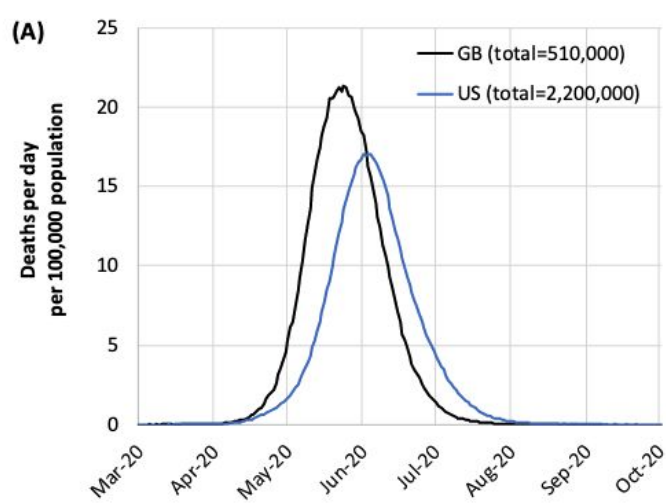


Generalizing Bulk-Synchronous Parallel Model for Data Science: From Data to Threads and Agent-Based Simulations

Zilu Tian¹, Peter Lindner¹, Christoph Koch¹, Markus Nissl¹, Val Tannen²
¹DATA, EPFL ²University of Pennsylvania

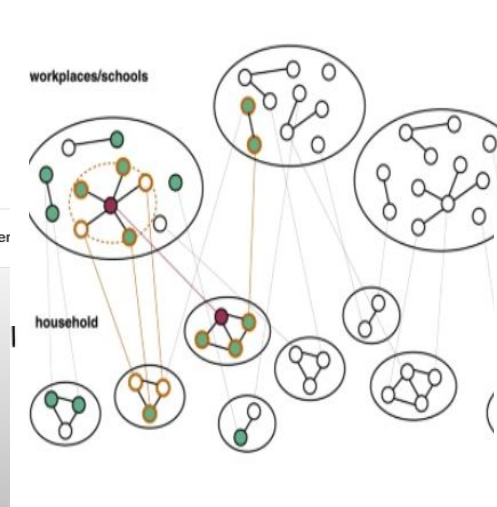
Why agent-based simulations matter?

• Epidemics



London's Imperial College predicts millions to die from coronavirus pandemic in UK and the US

Behind the Virus Report That Jarred the U.S. and the U.K. to Action



Special report: The simulations driving the world's response to COVID-19

Agent-based modelling of reactive vaccination of workplaces and schools against COVID-19

• Economics

Published: 05 August 2009

The economy needs agent-based modelling

J. Dooyne Farmer & Duncan Foley

Nature 460, 685–686 (2009) | Cite this article

13k Accesses | 661 Citations | 65 Altmetric | Metric

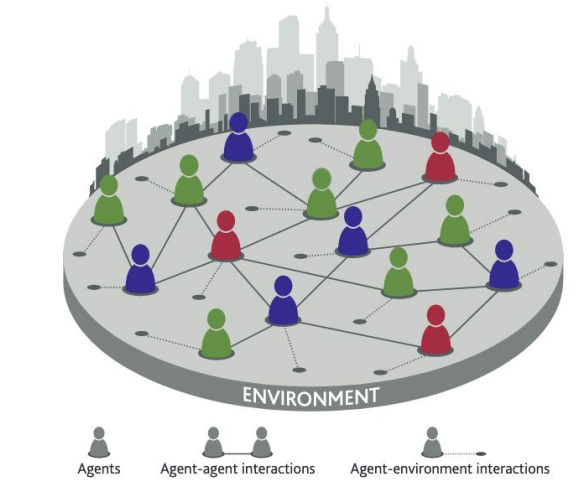


Chart 2 Reproducing stylised facts: the distribution of daily log-price returns

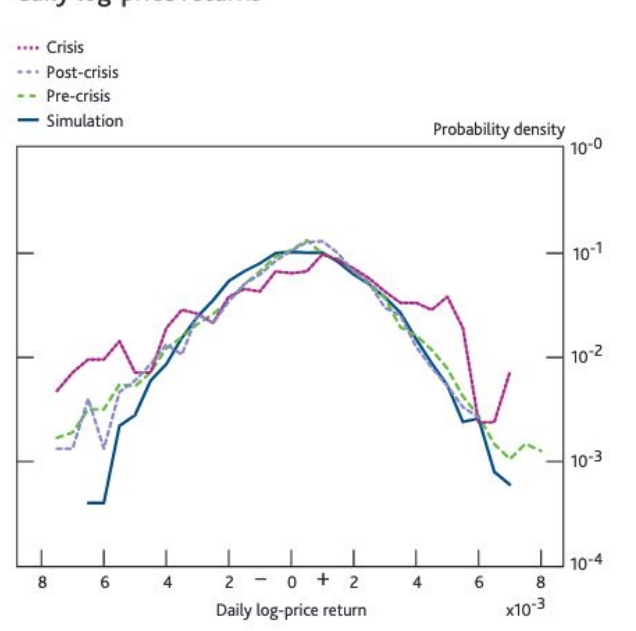


Chart 3 A benchmark run of the model showing boom and bust cycles in the house price index

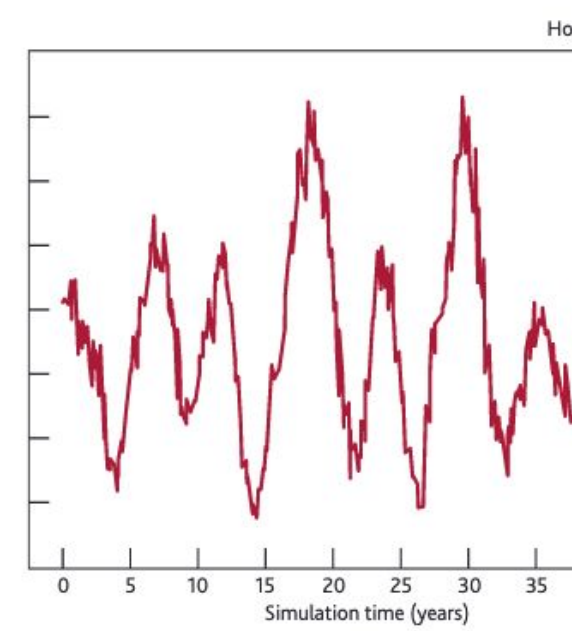
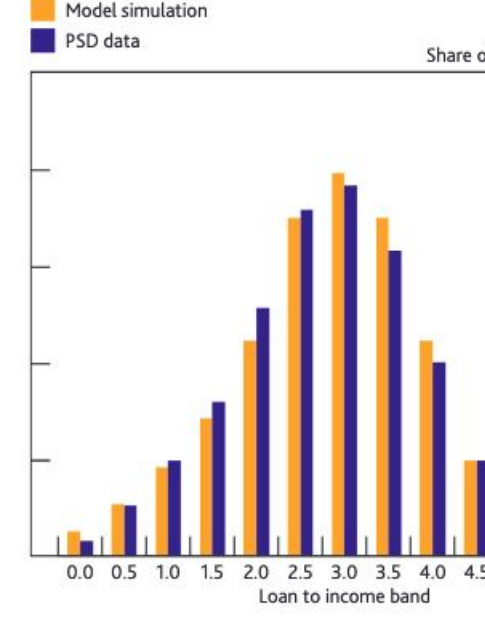
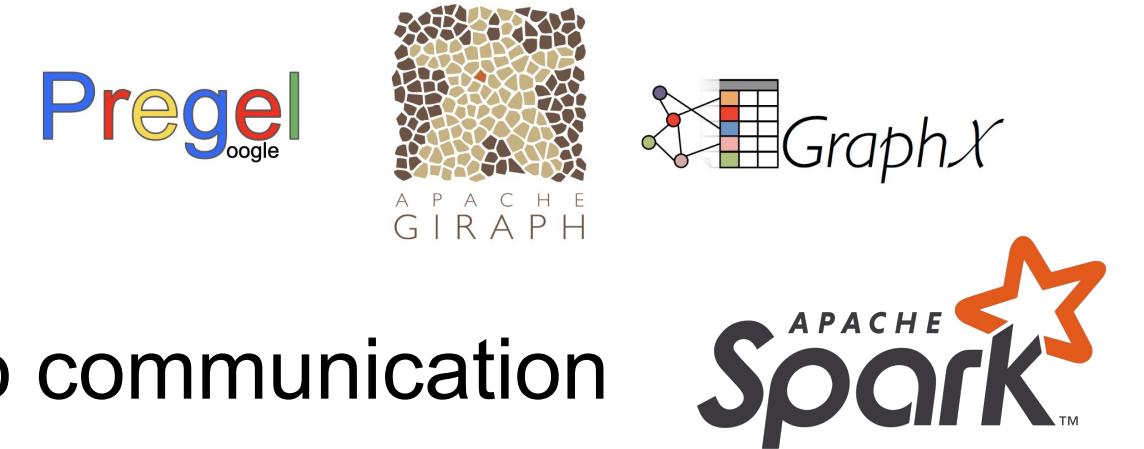


Chart 4 The model can reproduce the loan to income distribution for the United Kingdom



Challenges of agent-based simulations

- Agent-based simulations are flexible, but inefficient to execute
 - High concurrency
 - A realistic simulation has billions of agents
 - Code heterogeneity
 - “Think like a vertex” is homogeneous
 - Communication-intensive
 - Existing frameworks assume little or no communication
- For data management
 - Simulations generate a large amount of data
 - Image long-running simulations with billions of agents
 - Simulations form part of complex analytics pipelines
 - “How does the average wealth of the top 30% change?”
 - Simulations are can be viewed as model samples
 - “What is the average wealth of the population if we increase the initial wealth by 10%, 20% and 50%, respectively?”



What are agent-based simulations, really?

- Depend on who you ask!
 - A recent survey in 2020 listed
 - 36 general-purpose frameworks
 - 100+ specialized frameworks
 - Different assumptions about agents
 - NetLogo considers turtles as agents, along with patches and links
 - DMASON assumes each agent belongs to a temporal region
 - Repast Symphony assumes that agents actions are scheduled
 - Different assumptions about interaction
 - NetLogo assumes spatial-based interaction
 - DMASON is based on publish-subscribe paradigm
 - Repast Symphony allows instant changes to other agents' states

- The lack of formal models causes high heterogeneity
 - Increase users' learning curves
 - Decrease cross-platform result verification
 - Hard to select the right tool
 - Limit performance optimizations to framework-dependent

2 ABMS Software Packages				
2.1	AgentSheets			
2.2	NetLogo			
2.3	NetLogo			
2.4	NetLogo			
2.5	NetLogo			
2.6	NetLogo			
2.7	NetLogo			
2.8	NetLogo			
2.9	NetLogo			
2.10	NetLogo			
2.11	NetLogo			
2.12	NetLogo			
2.13	NetLogo			
2.14	NetLogo			
2.15	NetLogo			
2.16	NetLogo			
2.17	NetLogo			
2.18	NetLogo			
2.19	NetLogo			
2.20	NetLogo			
2.21	NetLogo			
2.22	NetLogo			
2.23	NetLogo			
2.24	NetLogo			
2.25	NetLogo			
2.26	NetLogo			
2.27	NetLogo			
2.28	NetLogo			
2.29	NetLogo			
2.30	NetLogo			
2.31	NetLogo			
2.32	NetLogo			
2.33	NetLogo			
2.34	NetLogo			
2.35	NetLogo			

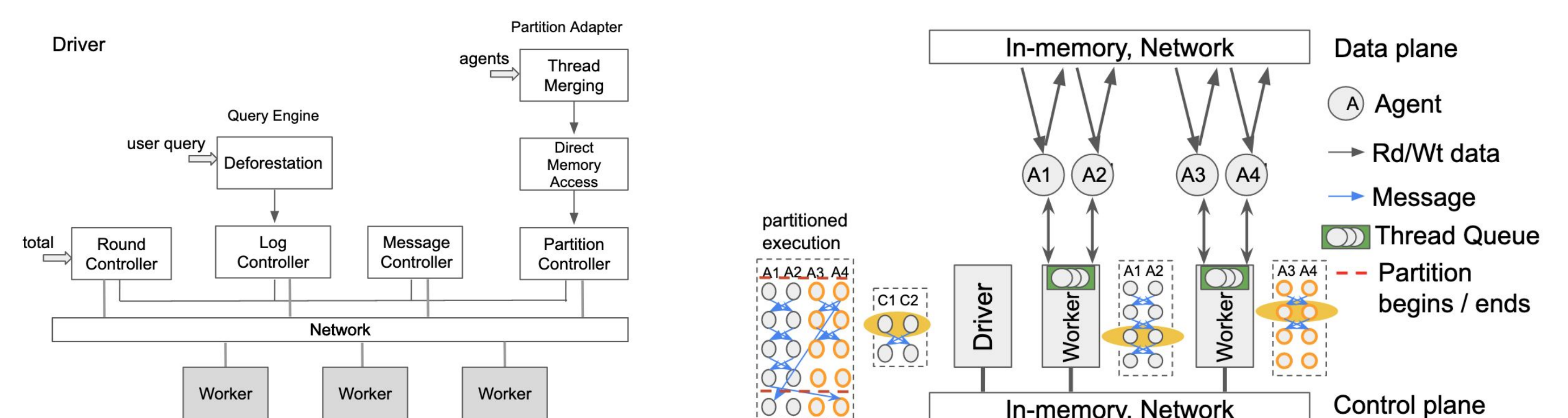
- Generally speaking, frameworks are round-based or asynchronous
 - whether agents proceed in lockstep
 - But frameworks have different flavors of “round-based” or “asynchronous”

Benchmark Description

- Population Dynamics
 - Simulate the game of life example in a 2D grid
 - Model each cell in the grid as an agent
- Economics
 - Simulate the bidding process in the stock market
 - Model traders and the stock market as agents
- Epidemics
 - Simulate individuals of states Susceptible, Infectious, Recovered, Hospitalized, or Deceased
 - Model the population and locations as agents
 - Use random graph models to simulate population connectivity
 - Erdos-Renyi Model (ERM)
 - Stochastic Block Model (SBM)

Contributions

- Formal models that define agents and their interactions
 - Programming model
 - Agents are sequential processes that communicate through messaging
 - A simulation is, conceptually, concurrent execution of interacting agents
 - "Simulate" as an operator for integrating with data science pipeline
 - Computational model
 - Weighted hierarchical BSP model
- Optimizations
 - Thread merging
 - Tame high-concurrency
 - Direct memory accesses
 - Bypass messaging overhead
 - Deforestation
 - Reduce the volume of generated data
- Implementation
 - An eDSL in Scala for parallel agent programming
 - A system architecture based on the BSP-like model



Performance

- Our system has on par or better performance than current BSP-like systems

