

From Automated Theorem Proving to Nuclear Structure Analysis with Self- Scheduled Task Parallelism

(a personal history of one programming model)

Rusty Lusk

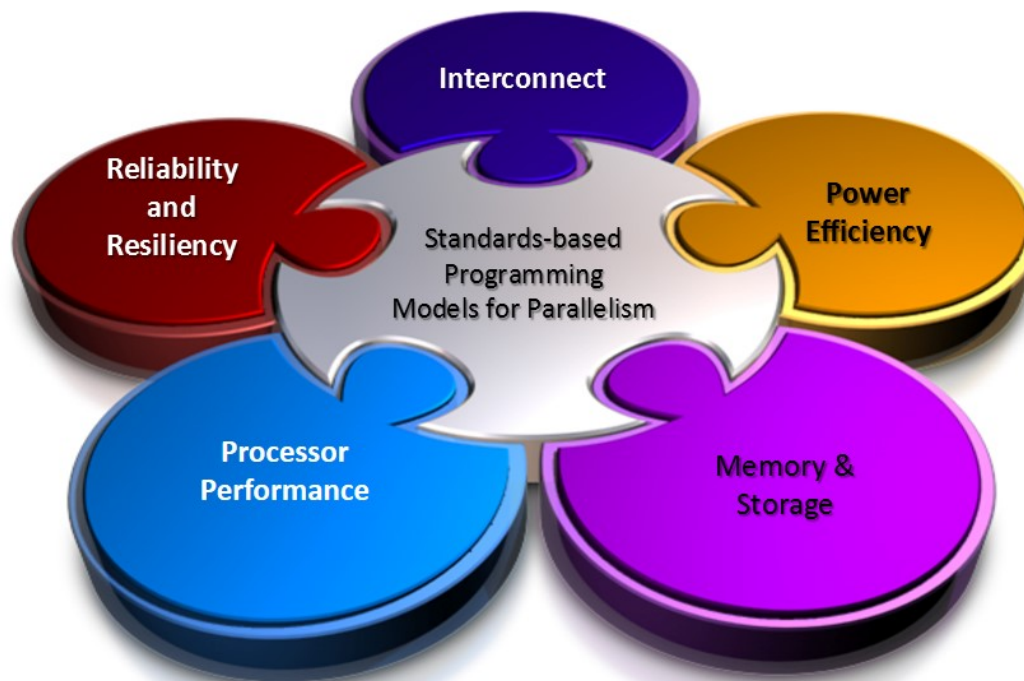
Mathematics and Computer Science Division

Argonne National Laboratory

My Theme

- *Application* programming models should be *simple*.
- Their *instantiations* might be more complex, and their actual *API's* might be even more so, not to mention their implementations.
- Also, the programming models of the libraries that implement them might be more complex.
- For example, the message passing *model* is simple: people are familiar with it from physical mail, phone calls, email, etc.
- There have been several *instantiations* (PVM, Express, EUI, p4, etc.) and multiple implementations of MPI.
 - As an application programming model, MPI is simple because applications use the simple parts
 - the more exotic parts of MPI are used by libraries to implement simple application programming models (or should be).
- MPI's full API is a really a system programming model, driven by library developers developing portable libraries that implement simple programming models for applications.





My Example

- I discuss here a simple programming model which has managed to remain simple through a number of instantiations and implementations.
- It is related to, but not the same as, several current task-based systems.
- It was how I wrote my first non-trivial parallel programs, back before the term “programming model” was in use (I didn’t know it was a programming model).
- I call it “self-scheduled task parallelism” (SSTP). My first work in computer science, after a stab at (very) pure mathematics, was in automated theorem-proving, at Argonne with Larry Wos, Ross Overbeek, and Bill McCune.
- The SSTP model was invented (not really on purpose) to parallelize the Argonne theorem prover (Otter).
- Therefore I am going to motivate it by entertaining you with a short introduction to automated theorem proving.



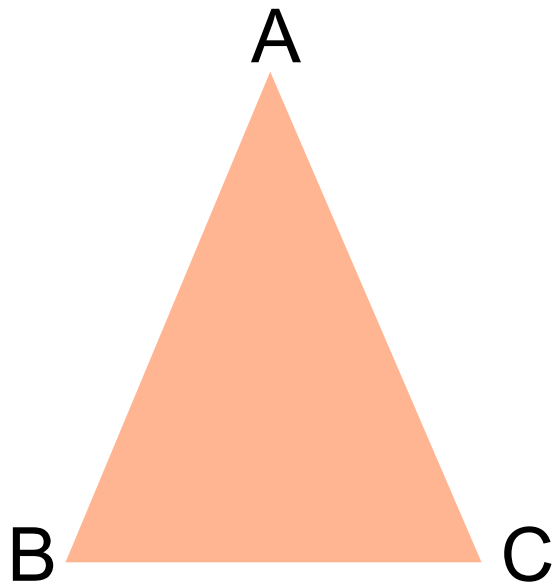
Outline

- Some ATP successes (why automated theorem proving is so much fun)
- Resolution-based automated theorem proving
 - How it works
 - A serial algorithm
 - Some parallel algorithms
- SSTP for a parallel Prolog system
- Why SSTP died out for a while
- Resurgence in Nuclear Physics SciDAC project as ADLB
 - Asynchronous Dynamic Load Balancing (ADLB), a minimal PM
 - ADLB is our current instantiation of the SSTP model
- Improving ADLB with another simple API, for memory management (DMEM)
- Recent results and current work

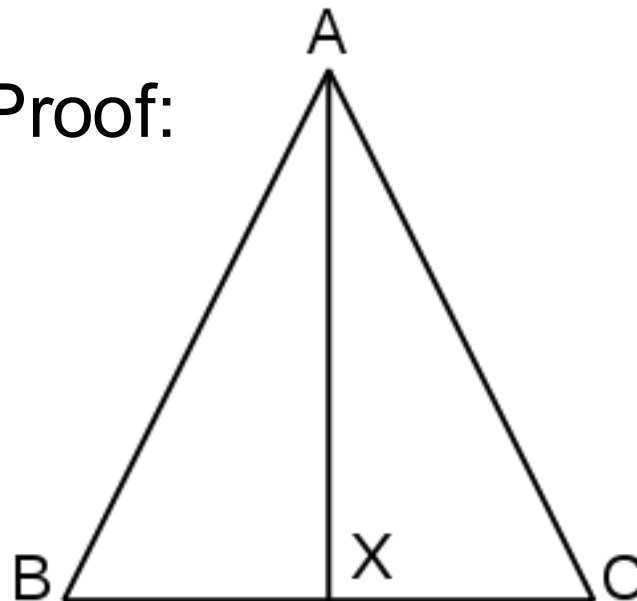


Going Way Back...

- Proposition 4 of Euclid's Elements (300 BCE) says that the base angles of an isosceles triangle are equal. This theorem is called the *Pons Asinorum**.



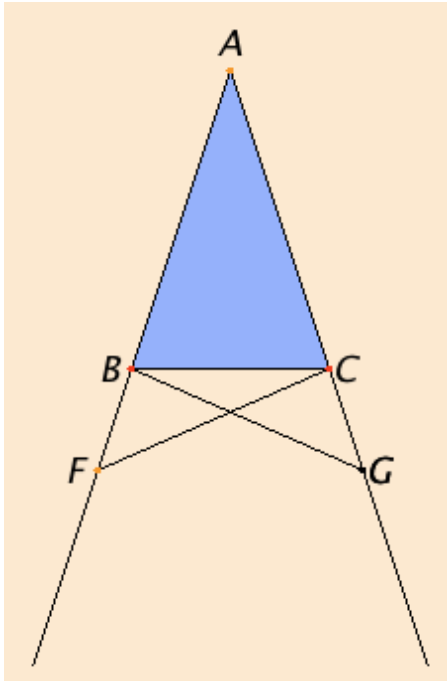
Proof:



Euclid

* "Bridge of Asses"

Euclid's Proof of the *Pons Asinorum* (From the *Elements*)

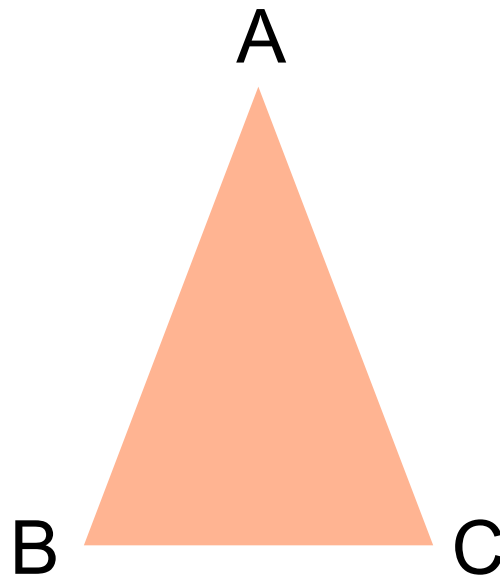


- Since AF equals AG , and AB equals AC , therefore the two sides FA and AC equal the two sides GA and AB , respectively, and they contain a common angle, the angle FAG .
- Therefore the base FC equals the base GB , the triangle AFC equals the triangle AGB , and the remaining angles equal the remaining angles respectively, namely those opposite the equal sides, that is, the angle ACF equals the angle ABG , and the angle AFC equals the angle AGB .
- Since the whole AF equals the whole AG , and in these AB equals AC , therefore the remainder BF equals the remainder CG .
- But FC was also proved equal to GB , therefore the two sides BF and FC equal the two sides CG and GB respectively, and the angle BFC equals the angle CGB , while the base BC is common to them. Therefore the triangle BFC also equals the triangle CGB , and the remaining angles equal the remaining angles respectively, namely those opposite the equal sides. Therefore the angle FBC equals the angle GCB , and the angle BCF equals the angle CBG .
- Accordingly, since the whole angle ABG was proved equal to the angle ACF , and in these the angle CBG equals the angle BCF , the remaining angle ABC equals the remaining angle ACB , and they are at the base of the triangle ABC . But the angle FBC was also proved equal to the angle GCB , and they are under the base.



A better proof, found by an automated theorem proving program in the 70's

- Triangle ABC is congruent to triangle ACB by the side-angle-side theorem. Corresponding angles of congruent triangles are equal. QED.



- Also Pappus, 320 CE



A More Recent Example

- The following open question was posed to our group by Irving Kaplansky, big-cheese algebraist at the University of Chicago:
- Is there a finite semigroup that has an anti-automorphism but no involution?
- Our program proved not only that answer was “yes,” but that the smallest was of order 7 and there were four such.
- Getting results publishable in math journals was even more fun than doing college algebra homework problems and theorem-proving benchmarks.



Kaplansky



How Resolution Theorem Proving Works

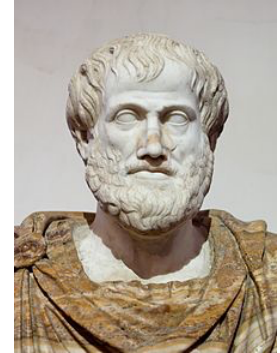
- Convert statements in 1st-order logic into disjunctive normal form, in which all variables are universally quantified and disjunction is the only connective. Implications become disjunctions:
- $\forall x, P(x) \rightarrow Q(x) \quad \longrightarrow \quad \neg P(x) \vee Q(x)$
- $\exists x, P(x) \quad \longrightarrow \quad P(a) \quad \text{(Skolem constant)}$
- Derive a new clause from 2 existing clauses by “cancellation”:

$$\begin{array}{c} \neg P \vee Q \\ P \\ \hline Q \end{array}$$

(Note: The original image has a large red 'X' over the entire derivation, indicating it is incorrect. The correct resolution rule would cancel P and $\neg P$ to yield Q .)



How It Works, continued



Aristotle

- Variables get instantiated to make the match:

All men are mortal.	$\neg \text{Man}(x) \vee \text{Mortal}(x)$
<u>Socrates is a man.</u>	<u>$\text{Man}(\text{Socrates})$</u>
Socrates is mortal	$\text{Mortal}(\text{Socrates})$

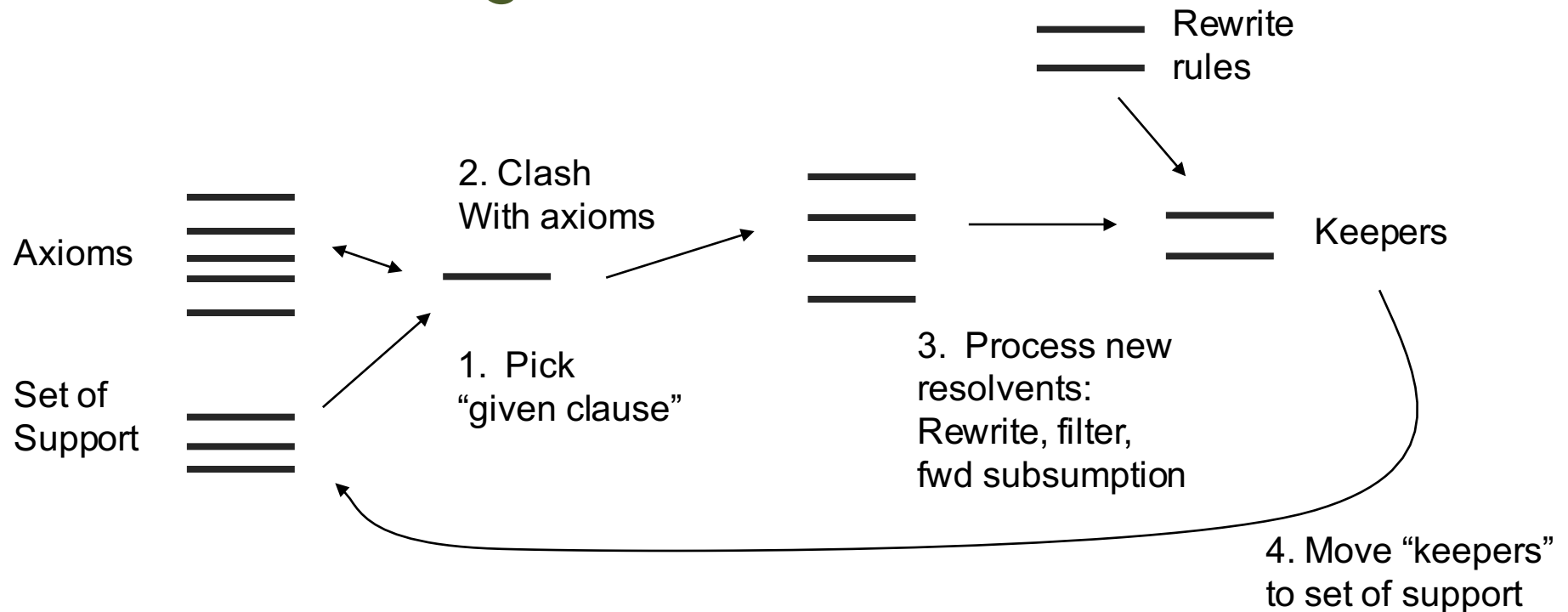
- To prove a theorem, state its denial and derive a contradiction, denoted by the “null clause.”

$$\begin{array}{c} P(a) \\ \neg P(a) \\ \hline \text{“ “} \end{array}$$

- The tricky bits are to avoid deducing too much and controlling redundancy



Otter's Basic Algorithm

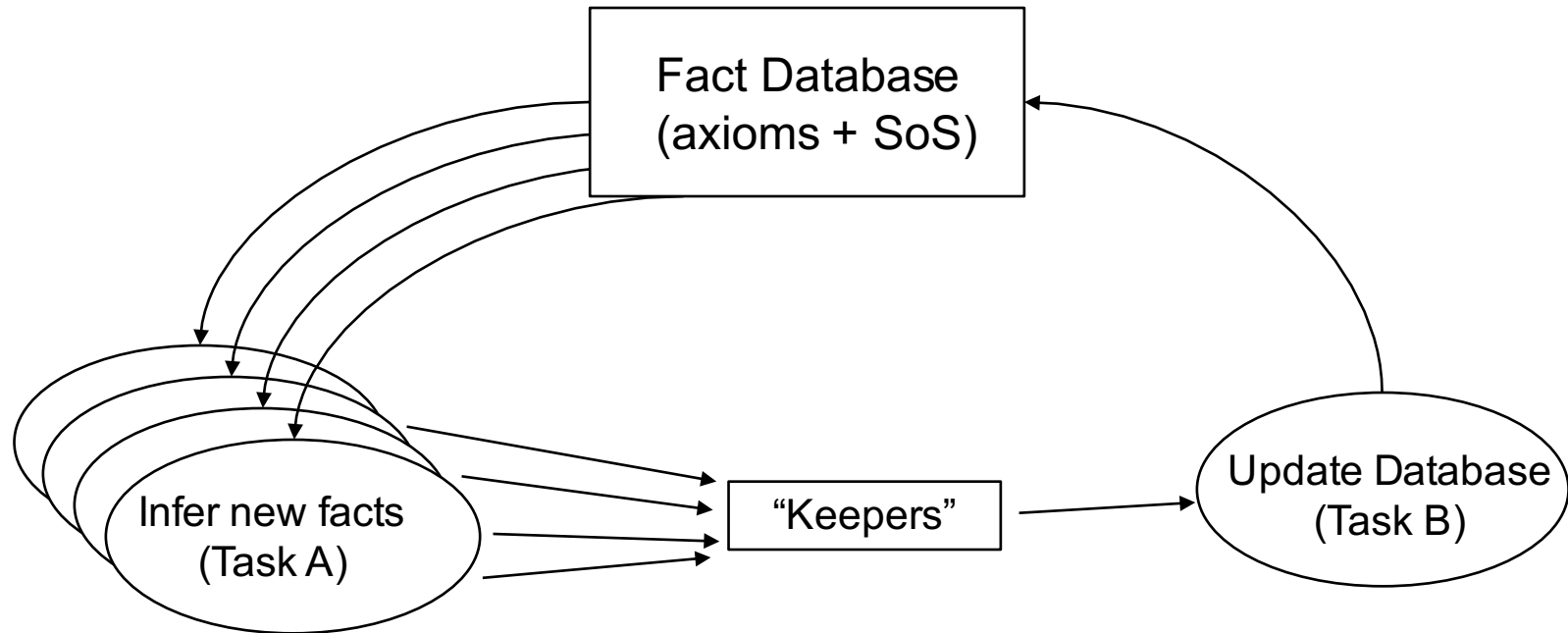


Repeat until you deduce the empty clause, SoS becomes empty, or you run out of time or memory.

- A *very* irregular computation
- First attempt at parallelism: process new resolvents in parallel
 - No good, since not enough parallelism, barrier before each new given clause
- Next version, process multiple given clauses in parallel



A Parallel Algorithm Without Deletion

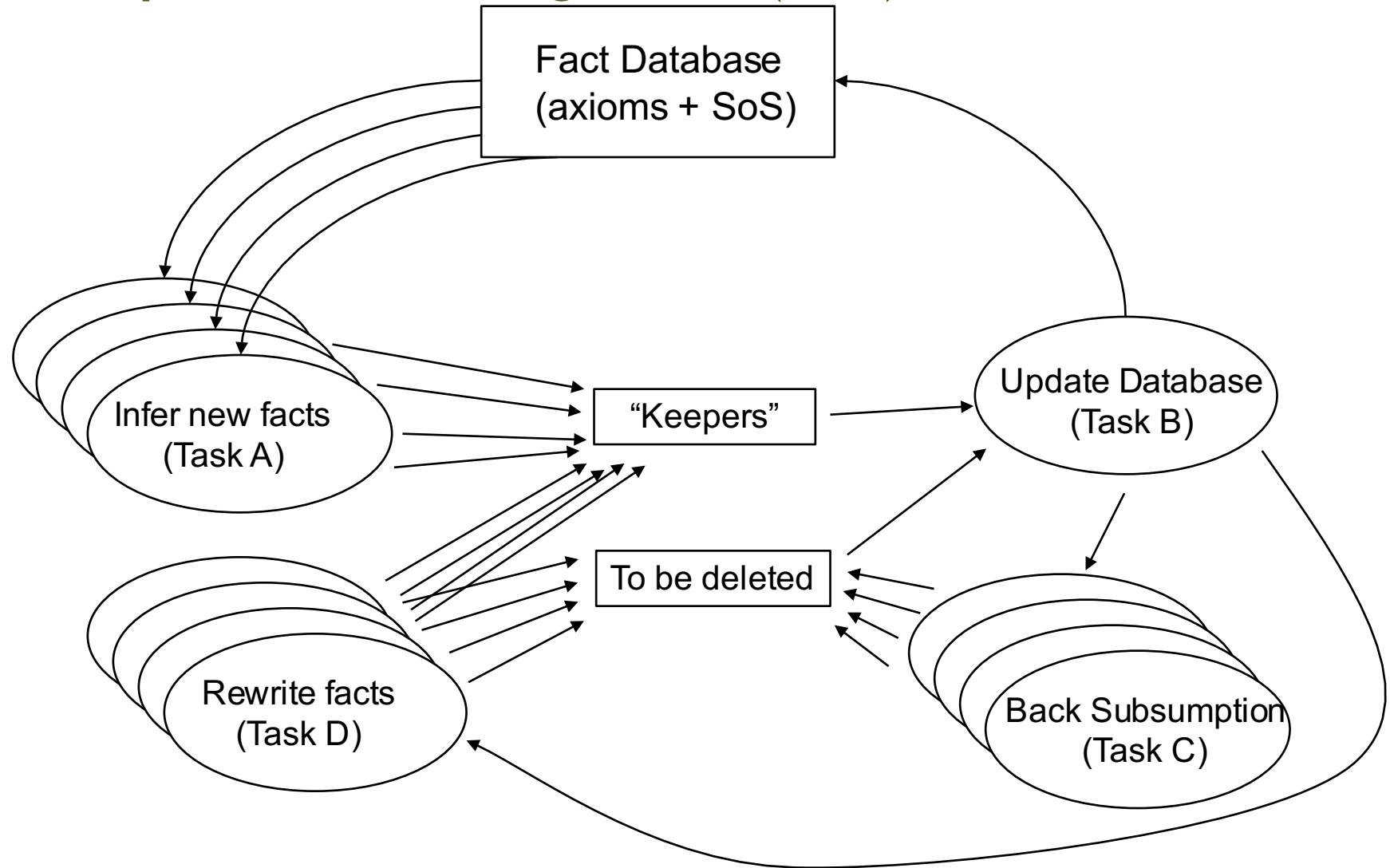


- Task A: Pick a given clause and carry out steps 1-3 from previous slide
- Task B: For each clause in K, do final forward subsumption test and add to set of support
- All processes:
 - If Keepers list is non-empty and no other process is doing Task B, do Task B
 - Else do Task A

This is the origin of SSTP.

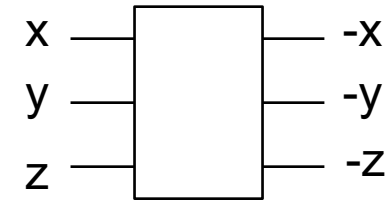


A Complete Parallel Algorithm (Roo)



- Only one process at a time does B, the rest is a free-for-all with no traffic cop or DAG
- Again, each process executes same loop, acquiring work, doing it, making new work

Some Old (But Good) Results



- The “two inverter” problem:
 - Design a circuit, using AND, OR, and just two NOT gates, whose 3 outputs are the inversions of its three inputs.
- In implicational propositional calculus, the law of hypothetical syllogism can be derived from a proposed single axiom by condensed detachment:
 - $\neg P(x) \vee \neg P(i(x,y)) \vee P(y)$ (Condensed detachment)
 - $P(i(i(i(x,y),z), i(i(z,x), i(u,x)))$ (Lukasiewicz axiom)
 - $\neg P(i(i(a,b), i(i(b,c), i(a,c))))$ (Denial of hypothetical syllogism)

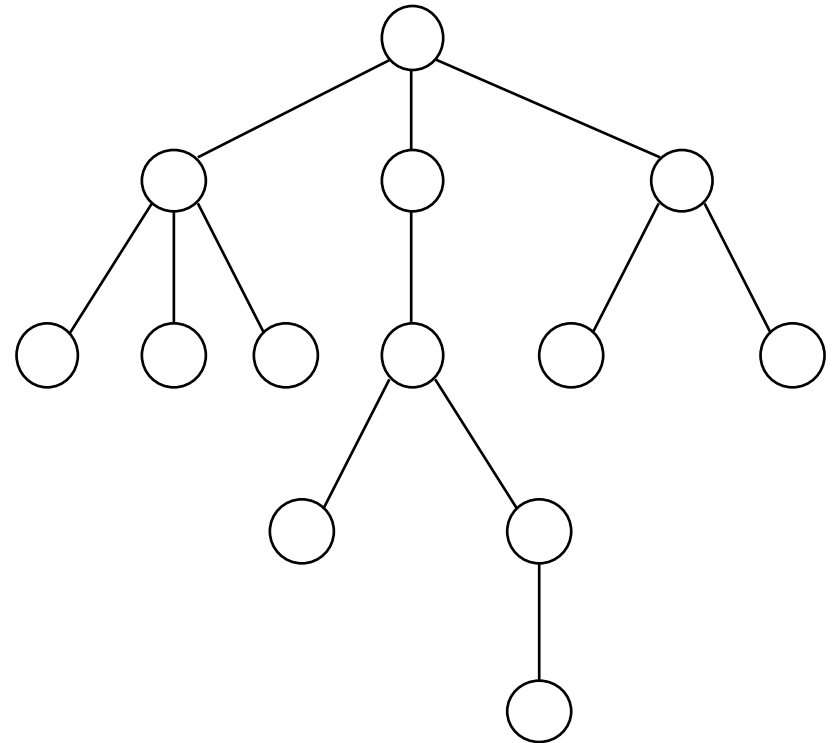
	Otter-2inv	Roo24-2inv	Otter-Luka	Roo24-Luka
Runtime (sec.)	47236	2237	29098	1269
Generated	6323644	6351410	6706380	7108289
Kept	21342	21343	20410	18759
Speedup	1.0	21.1	1.0	22.9



- On 24-processor Sequent

Parallel Prolog

- Creating/acquiring work is again done by modifying a shared data structure
- Just beginning to identify and abstract these operations into general putting work into, and getting work out of, a shared work pool



SSTP Takes a Vacation

- As the number of processors multiplied, shared memory couldn't scale, and large-scale parallel computing went to message passing.
- DOE lost interest in inference as the hope of a program verification miracle faded.
- SSTP evolved (backwards) into the manager-worker programming model (e.g. Linda)
- This solved beautifully the load-balancing problem for irregular computations but hit its own scalability problem
 - Too many workers for a single manager to keep up with
 - Too little memory for a single manager to store the structures defining the work pool



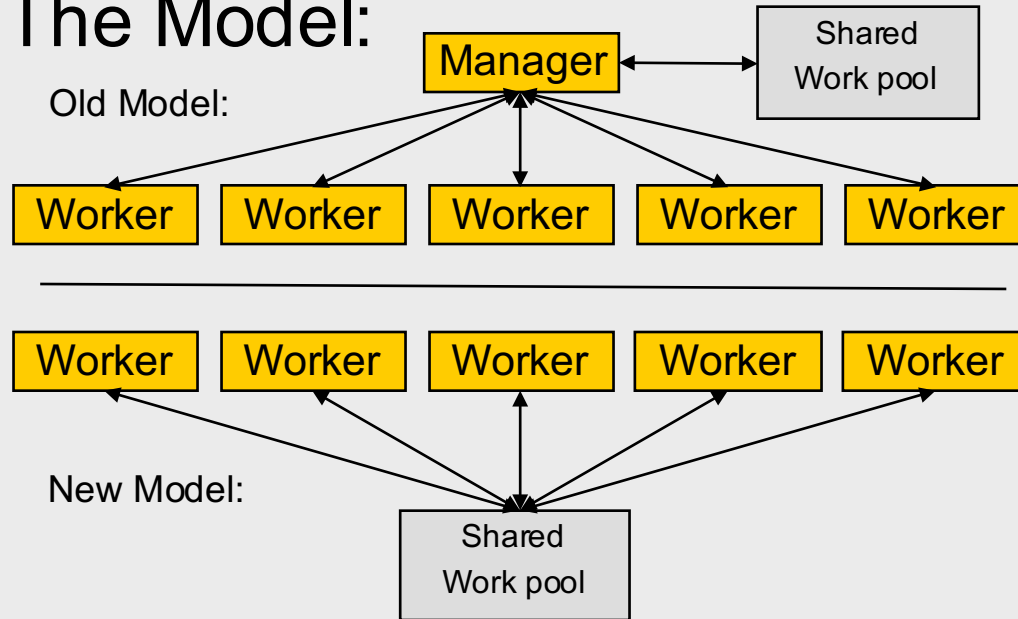
Green's Function Monte Carlo - A Complex Application

- Green's Function Monte Carlo -- the “gold standard” for *ab initio* calculations in nuclear physics at Argonne (Steve Pieper, Physics Division)
- A non-trivial manager/worker algorithm, with assorted work types and priorities; multiple processes create work dynamically; large work units
- Had scaled to 2000 processors on BG/L, then hit scalability wall.
- Needed to get to 10's of thousands of processors at least, in order to carry out calculations on ^{12}C , an explicit goal of the UNEDF SciDAC project.
- The algorithm threatened to become even more complex, with more types and dependencies among work units, together with smaller work units. An extremely irregular computation.
- Wanted to maintain original manager/worker structure of physics code
- This situation brought forth the Asynchronous Dynamic Load Balancing Library (ADLB), giving up generality for scalability and ease of use.
- Achieving scalability has been a multi-step process
 - balancing processing
 - balancing memory
 - balancing communication
- Now runs with 100's of thousands of processes



ADLB On One Slide

The Model:

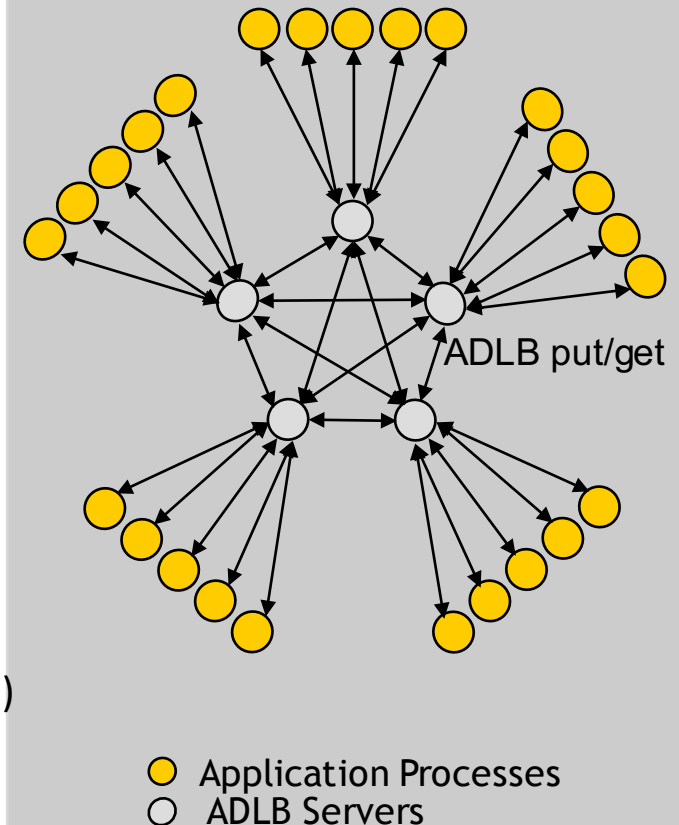


The API:

- ADLB_Put(type, priority, len, buf, target_rank, answer_dest)
- ADLB_Reserve(req_types, handle, len, type, prio, answer_dest)
- ADLB_Get_Reserved(handle, buffer)
- and a few housekeeping calls...

ADLB abstracts the idea of creating/acquiring work using put/get of work units into a work pool

An Implementation:



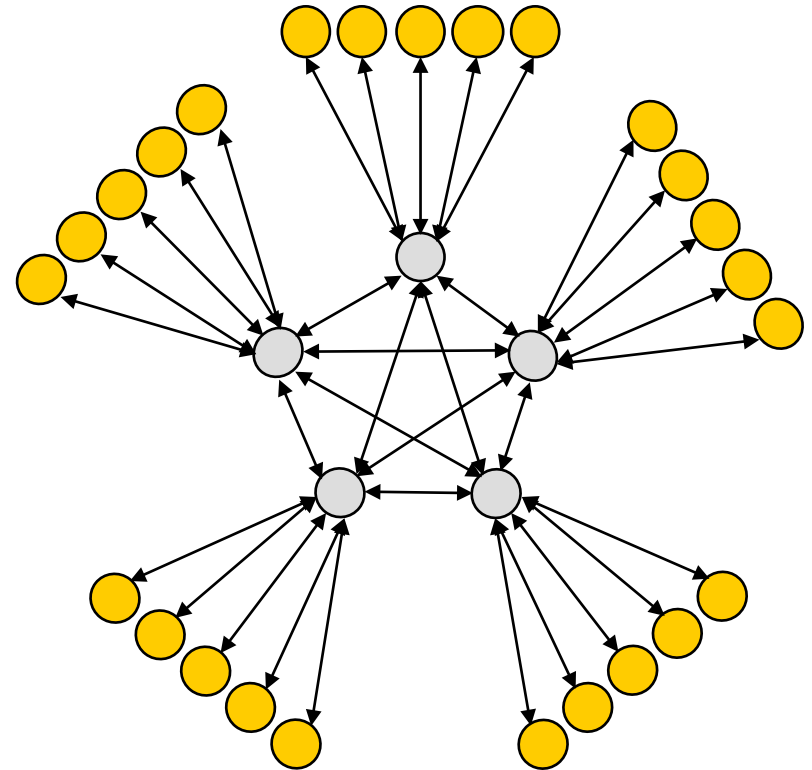
ADLB Uses Multiple MPI Features

- ADLB_Init returns separate application communicator, so application processes can communicate with one another using MPI as well as by using ADLB features.
- Servers are in MPI_Iprobe loop for responsiveness.
- MPI_Datatypes for some complex, structured messages (status)
- Servers use nonblocking sends and receives, maintain queue of active MPI_Request objects.
- Queue is traversed and each request kicked with MPI_Test each time through loop; could use MPI_Testany. No MPI_Wait.
- Client side uses MPI_Ssend to implement ADLB_Put in order to conserve memory on servers, MPI_Send for other actions.
- Servers respond to requests with MPI_Rsend since MPI_Irecv are known to be posted by clients before requests.
- MPI provides portability: laptop, Linux cluster, BG/Q, Cray
- MPI profiling library is used to understand application/ADLB behavior.



© 2013 Pearson Education, Inc. or its affiliate(s). All rights reserved. Pearson Education, Inc., publishing as Pearson Benjamin Cummings, 101 Philip Drive, Assinippi Park, New York, NY 10984-2135. All rights reserved. Printed in the United States of America. This publication is protected by copyright. Permission is granted to reproduce copies of this publication for personal or internal use, on the basis of payment of \$1.00 per copy to the Copyright Clearance Center, Inc., 222 Rosewood Drive, Danvers, MA 01923. For those organizations that have been granted a photocopy licence by CCC, a separate system of payment has been arranged. The fee code for users of the CCC Transactional Reporting Service is 0893-4003/13 \$1.00. ISBN 0-321-83099-5. Printed on acid-free paper.

- The multiple servers were originally introduced to spread the communication (and computational) load that were swamping the one master.
- But they also store the data for the work units.
- As the work units became larger, we needed more servers for their storage capability, exacerbating the synchronization problem.
- Solution: decouple work unit allocation from work unit storage.



DMEM - A library to provide a shared-memory model on a distributed-memory machine

- API summary: put, get, copy, free, get-part, update
- User (application or another library) refers to a memory object via a (small) *handle*, which encodes its location and size.
- DMEM runs as a separate thread in applications, sharing memory with application processes, so local operations are fast.
- Optimization: put and copy operations are local if possible.
 - For non-local operations, multiple optimization strategies are possible
- Looking ahead, object size is of type MPI_Aint, which is typically a long int in C and an integer*8 in Fortran.



DMEM API and Implementation

- The API:
 - `DMEM_Put(void *, MPI_Aint, dmem_handle);`
 - `DMEM_Get(dmem_handle, void *);`
 - `DMEM_Get_part(dmem_handle, MPI_Aint, MPI_Aint, void *);`
 - `DMEM_Copy(dmem_handle, dmem_handle);`
 - `DMEM_Update(dmem_handle, MPI_Aint, MPI_Aint, void *);`
 - `DMEM_Free(dmem_handle);`
- The Implementation:
 - At `DMEM_Init`, each process forks a thread, which processes DMEM requests independently of the application thread.
 - There is a manager who keeps an approximate view of memory used by DMEM on each process.
 - So far, the manager has not become a bottleneck.
 - Optimization problem: on `DMEM_Put`, where should the memory be allocated? (locally if possible, then round-robin, load balanced, or mixture)
 - Relies on MPI for communication



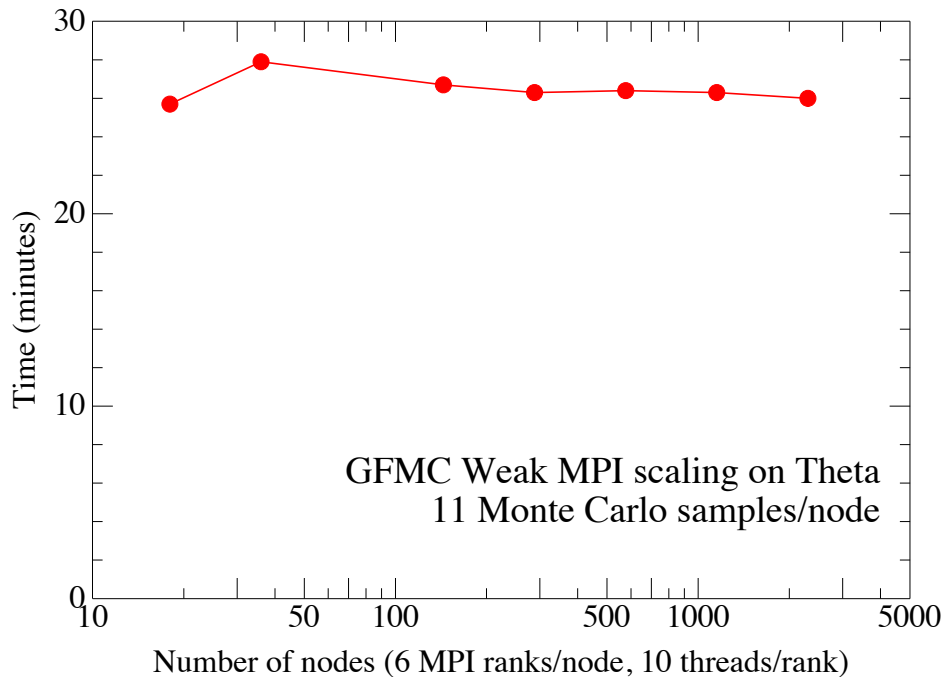
How DMEM Helps ADLB

- DMEM's MPI communicator is all of GFMC's client (application) processes (not the servers).
- GFMC is modified to store work units containing DMEM handles instead of the large blocks of data that used to be the work units; data is stored and retrieved via DMEM_Put and DMEM_Get.
- All of application processes' total memory is now available.
- Work units presented to ADLB are tiny (contain handles instead of entire work unit data).
- So way fewer servers are needed for storage.
- So ADLB's synchronization challenge disappears.
- Everybody wins!



Recent Results for GFMC with DMEM on Theta

- Theta is new Knights Landing – based machine at Argonne, with 64 cores/node, and we have just started experiments
- GFMC is hybrid: OMP + ADLB + DMEM
 - Strong OMP scaling per node up to # cores (1 MPI rank on node)
 - Better throughput with multiple ranks per node
 - Weak scaling with ADLB up to current size of machine



Each rank uses 9
OMP threads and
one pthread for
DMEM; 6 ranks
per node



A Lurking Future Problem (LFP)

- (Near) future machines are going to have lots of memory per node (for huge work units) and lots of threads (hardware and software) per node (to work on them).
- What if an ADLB (or even just a DMEM) application wants to utilize work units whose size is larger than 2 GB (approximately the size of a 32-bit integer)?
- ADLB and DMEM are agnostic about the internal structure of work units, so their internal messages use MPI_BYTE as their message type, so the count argument in MPI communications is the size (in bytes) of the message.
- MPI_{Send/Recv} specifies the count argument as an integer (still 32 bytes on most systems).
- The MPI-3 forum decided not to change this, because “long” messages could be sent/received on an MPI-compliant implementation by using MPI datatypes to lower the count argument into the 32-bit range.
- But:
 - Some people (even me, an MPI enthusiast) consider MPI datatypes inconvenient.
 - Some important MPI implementations are not MPI-compliant! (e.g. Mira and Titan)
- Solution: a long-message library for anyone who needs it: MPIL
 - Looks like MPI, except for MPIL_Count in MPI_Send/Recv, etc.
 - Limited version (enough for DMEM) working now



MPIL - MPI Long Messages

- API
 - MPIL_Init(comm)
 - MPIL_Send(*buf, MPI_Count count, datatype, rank, tag, comm)
 - MPIL_Recv(*buf, MPI_Count count, datatype, rank, tag, comm, MPIL_Status &status)
 - MPIL_Finalize(comm)
 - MPIL_Probe(...) (the tricky one)
 - MPIL_Bcast(...) (etc.)
- Implementation (in progress)
 - For MPI-standard-conforming implementations:
 - Construct datatype consisting of large number of user's datatypes
 - MPI_Send/Recv using this datatype and 32-bit value of count. Use (hidden) struct datatype if division has remainder
 - For implementations where the underlying communication layer can only handle 32-bit-size messages:
 - Divide user message into multiple smaller messages (chunks).
 - Send header with first chunk, so MPIL_Recv knows how many MPI_Revs to post.
 - Use hidden communicator to help with MPIL_Probe.



Summary

- Automated proving, an irregular computation, motivated our initial self-scheduling, load-balancing approach
- ADLB, its current instantiation, demonstrates that by giving up some generality, a programming model can provide scalability without complexity for (some) applications.
- GFMC motivated ADLB, which motivated DMEM, which motivated MPIL.
 - But all 3 are small, portable, independent libraries
- DMEM was a big help to ADLB, but is potentially useful in a more general context. (e.g. to exploit multiple types of memory in a hierarchical memory system). Needs wider user input.
- MPIL will be a simple, portable way to provide long message support to any MPI program at lowest cost.
- Even little-bitty libraries (i.e. with small API's) can be useful in HPC physics applications (as long as they have a Fortran interface, of course).
- Automated theorem proving might be currently somewhat out of fashion, but wouldn't it be great if we could....

Make America logical again!



The End

