

climateprediction.net and BOINC



Tolu Aina¹
Carl Christensen^{1,2}
Neil Massey²



University of Oxford

¹Department of Atmospheric, Oceanic, and Planetary Physics

²Oxford University Computing Laboratory



climateprediction.net

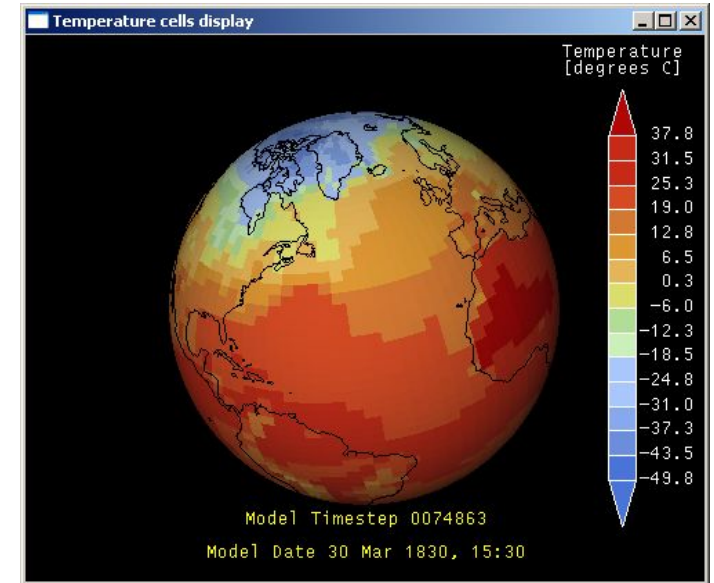
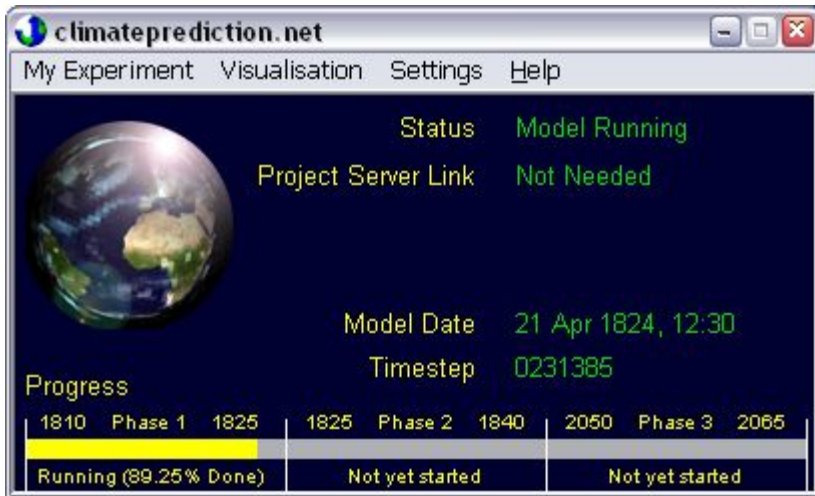
The Challenges...

- Climate models (ESM's, AOGCM's etc) are very large, complex systems developed by physicists sometimes over decades (& proprietary in case of UKMO)
- ~1 million lines of Fortran code (HadSM3 -- 550 files, 40MB text source code)
- Little documentation (the science is well documented by not the software and design of the system per se)
- Also “utility” code written by various scientists & students over the years (outside of model code, 220 files, 12MB source, 250K lines); often workable but hard to implement on a cross-platform PC project
- Meant to be run on supercomputers, primarily 64-bit – not designed (or indeed envisioned) to be run on anything other than a supercomputer or at the very least, a Linux cluster

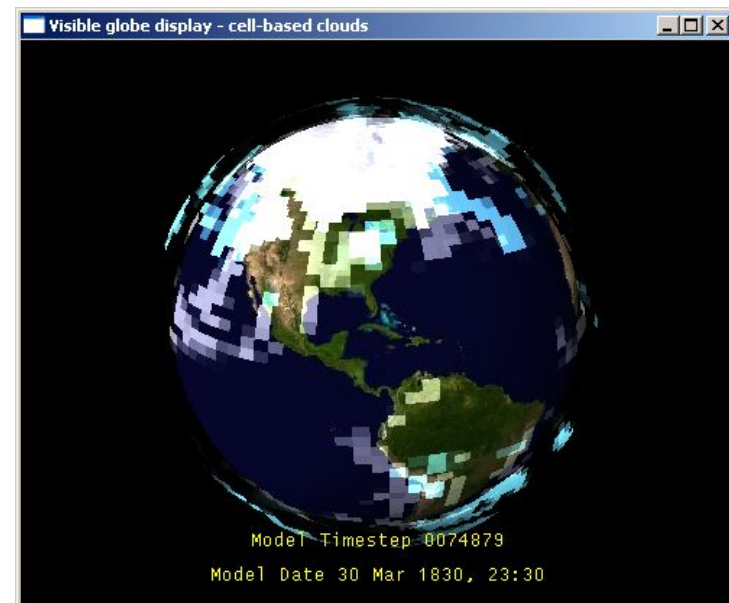
The Solution (“pre-BOINC”)

- CPDN spent about a year just getting the model to run on a Windows PC, and verify results were accurate compared to supercomputer & cluster runs etc
- With this “proof of concept” CPDN got funding to make the full-fledged project
- ...then spent another 1 - 1.5 years making a “vertical application” to distribute model runs, stats, credits...
- This was “launched” at an event at the Science Museum in London on 12 September, 2003
- “classic CPDN” 3 million model-years to date; akin to having full use of the Japanese Earth Simulator since the launch

“Classic” Client & Server Layout



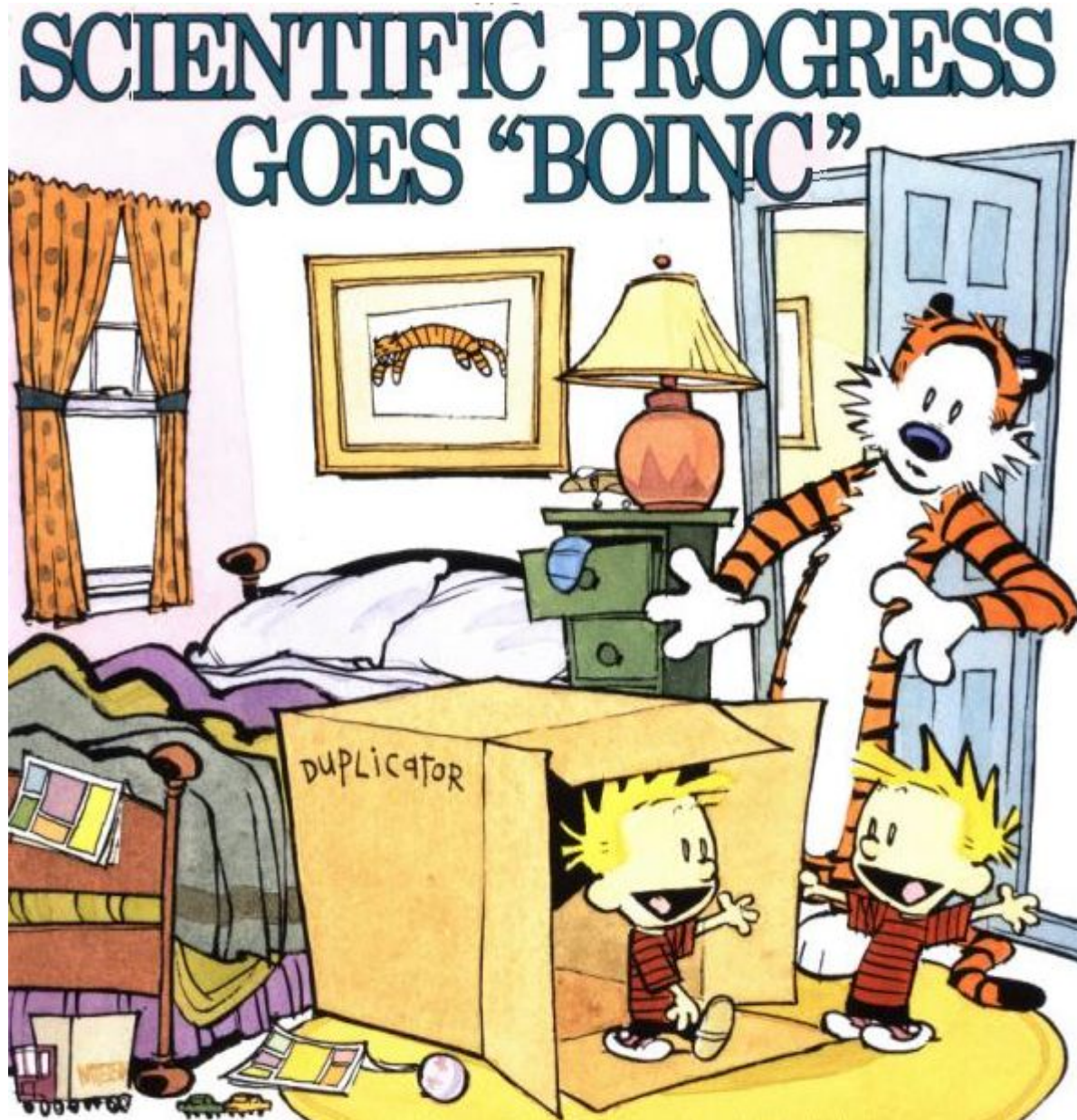
- GUI developed in MS Visual C++
- GUI controls the Model which is Compaq Visual Fortran (ported from MetOffice model ~million lines of Code)
- Server is Linux (RedHat9) with Oracle 9i as the database
- Client and Server together represent a custom “vertical application” for distributed computing



Lessons Learned From “Classic” CPDN

- User attrition rate is high
 - out of 47,000 total participants active “trickling” machines is about 17,500 after 3 months since the public launch of 12th Sept '03, and stabilized at 10K before the BOINC port
- The model should be easy to run, but it will always be a disk & memory & CPU “hog”
 - 50MB RAM usage, 1GB hard disk, CPU usage is at low “idle” priority, disk I/O about 5-10GB/day.
- We rushed to get user pages done since we were “bogged down” in getting a stable model/client out the door
- We possibly underestimated the enthusiasm of people for “raw stats” since we were focused more on the scientific outcomes.
- Some things we have no control over – SETI & Prime & Folding & other DC projects were multi-platform and able to run on “low power” PC’s – a climate model is much harder to port & split up into “DC-sized” components, doesn't run reliably < Pentium III.

CPDN and BOINC Integration

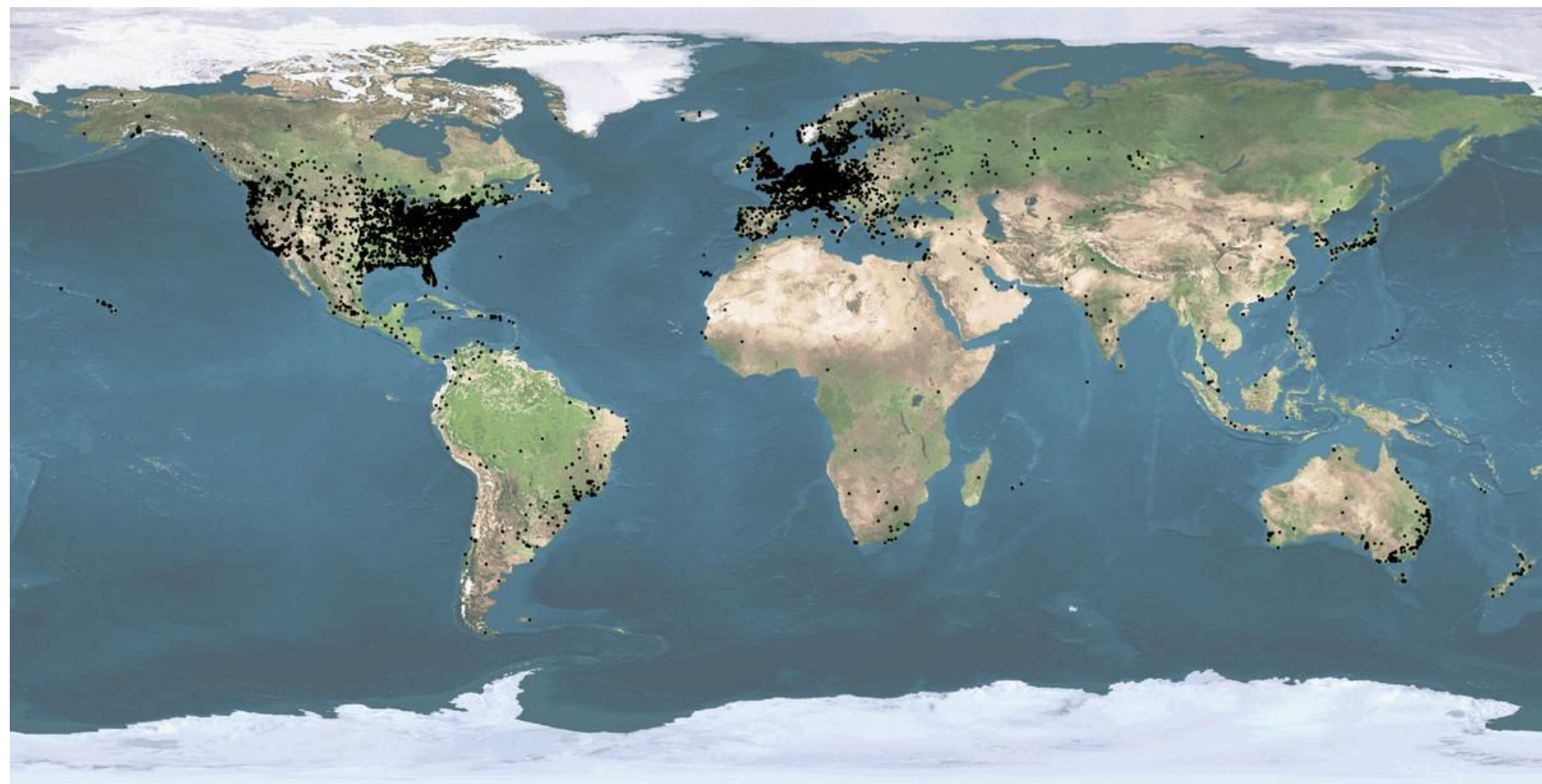


*Apologies to
Bill Watterson*

The Solution - BOINC!

- Known limitations on our “vertical app” - only Windows, geared only towards the one model (HadSM3), etc
- Visited Berkeley in December of '03 – state of BOINC was still early for production, but looked promising
- Tolu & Carl spent first 6 months of '04 porting HadSM3 to BOINC – including port to Linux & Mac
- Compare with our vertical app that took 1.5 years (on top of the time to get the model to run in Win32)
- To date 54K completed runs (4.2mn model-years) under BOINC, 50% more than original CPDN
- All future development to be under BOINC (sulphur-cycle, coupled model, HadAM3, MIT/CCSM)
- BOINC allows us to spend our time on the climate science & modelling issues, not so much on the “infrastructure” -- better use of time for CPDN staff

***climateprediction.net* BOINC Users Worldwide**
80,000 users total: 30,000 at any one time (trickling)



As CPDN Principal Investigator Myles Allen likes to say...
“this is the world's largest climate modelling supercomputer”

CPDN / BOINC Issues

- CPDN is a “compound app” in BOINC -- a very low-CPU load “monitoring” thread (C++) that tracks the high-CPU load climate model (primarily Fortran)
- Models sensitive and can crash anytime, anywhere!
- Run time of HadSM3 on a 2GHz P4 Win32 machine is about 4 weeks – far longer than any other app
- So “original CPDN” had trickles to give basic stats/info to our servers (to see “who's still alive?”)
- This (& compound app support) was implemented in BOINC in version 4 and is critical for CPDN
- We also use zip files to distribute app & workunits
- CPDN “state files” are very large & numerous, so procs run out of the projects/cpdn directory, not slots/* dirs

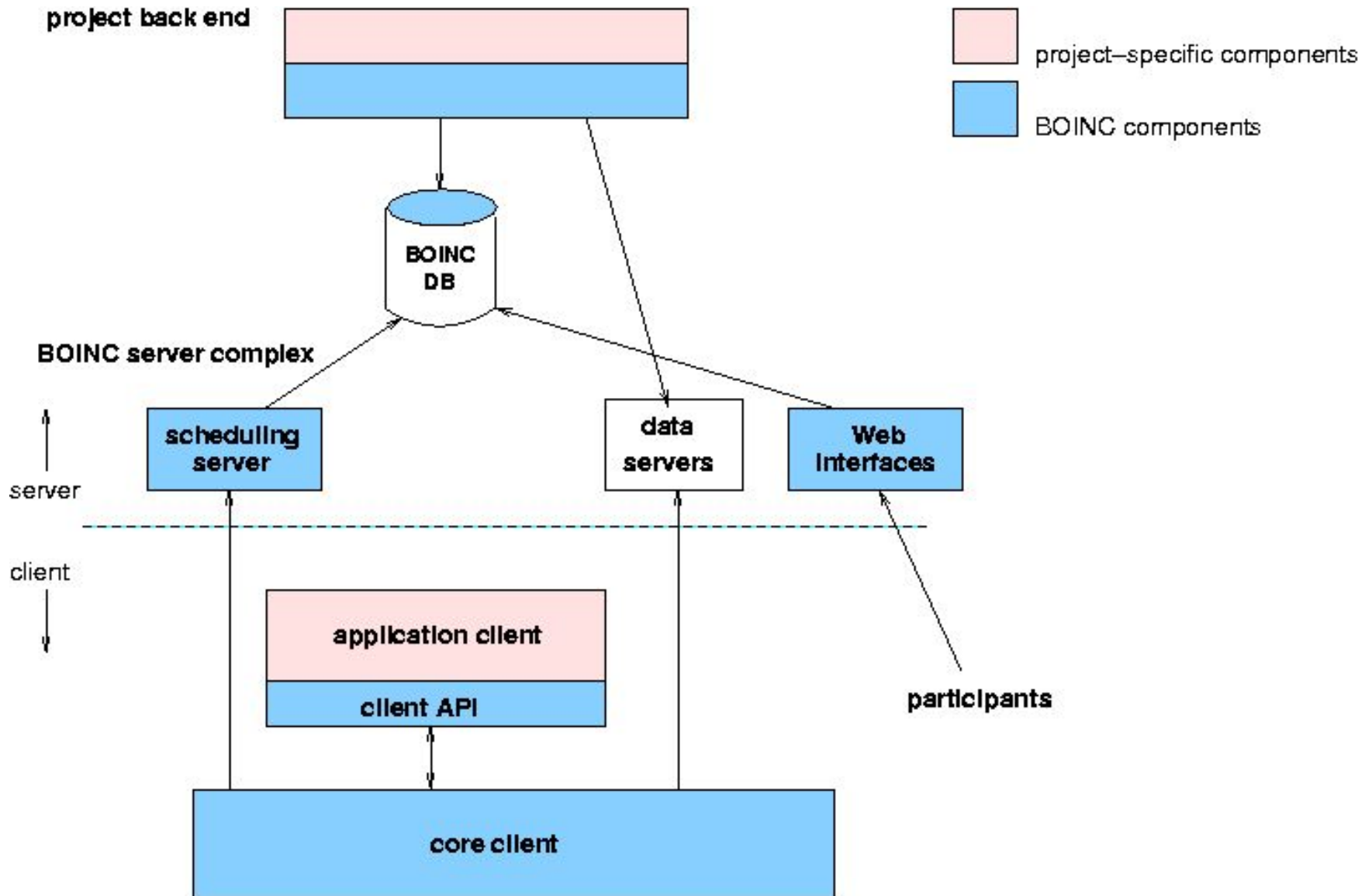
Why BOINC? (or “I Dream of SETI”)

SETI@Home was the pioneer and “killer” DC app with user & CPU-year numbers that are the envy of every other DC app, 5 million total users, 500K at once

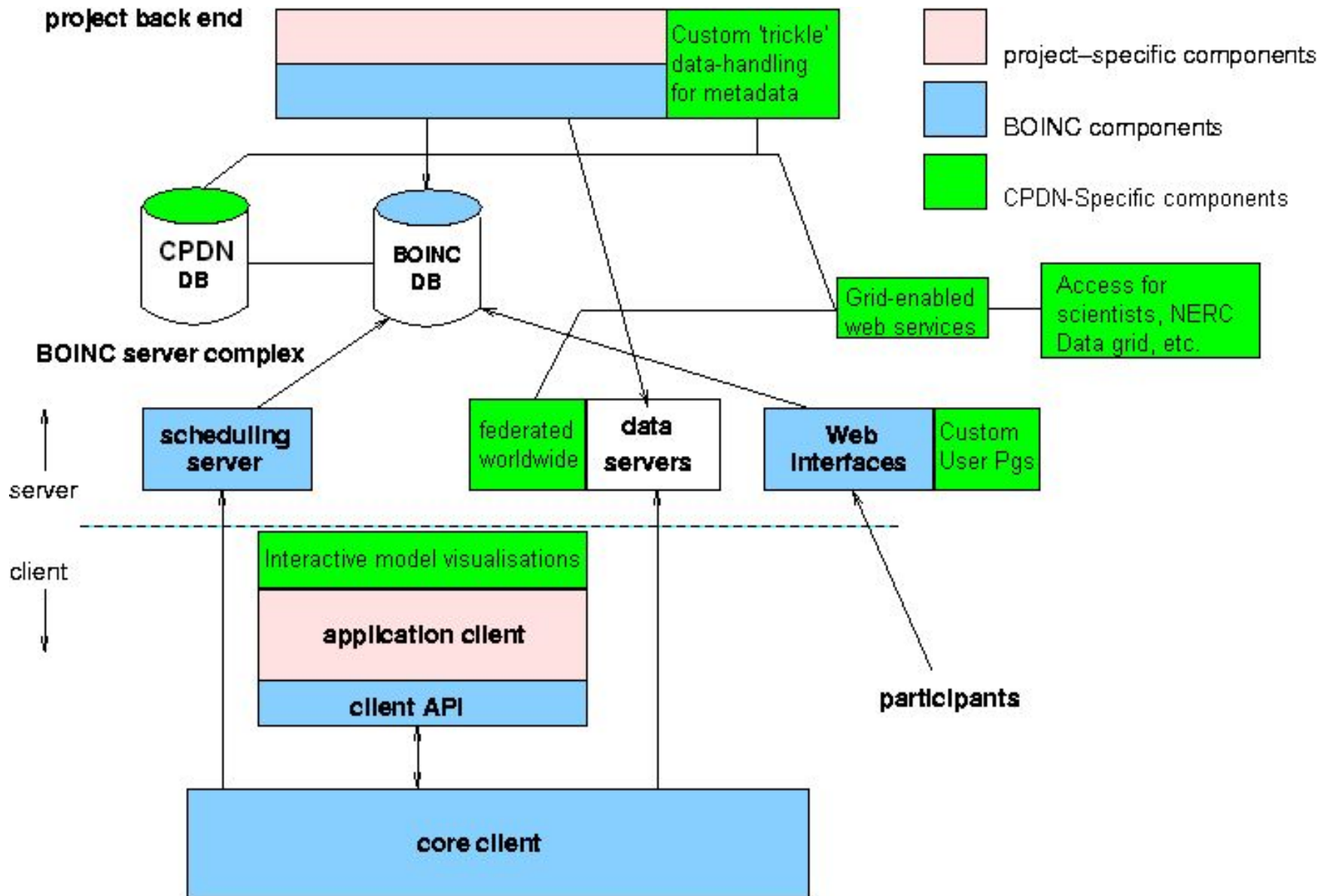
- BOINC is based on the experiences of the SETI@Home team in handling millions of users, downloads and uploads, investment of >US\$1million
- So it makes sense to use BOINC which has a “tried and tested” framework instead of keep playing “catch-up” and “reinventing the wheel”
- Basically, BOINC will allow us to focus on what we do best (or should be doing best):
 - Climate science, climate modelling, visualisation packages (peer-to-peer perhaps?), cross-platform porting of models, grid applications to “clamp on” the BOINC server-side
- Reality: Everybody wants “old SETI #'s” of millions of users – but tens of thousands more likely (but “nothing to sneeze at” -- it's still 'Japanese Earth Sim' level!)



BOINC Roadmap



BOINC Roadmap with CPDN



CPDN / BOINC Server Setup

- Database Server – Dell PowerEdge 6850, two Xeon 2.4GHz CPUs, 3GB RAM, 70GB SCSI RAID10 array (RAID5 originally on this former Oracle server, but seems sluggish for MySQL).
 - my.cnf configured to use about all of the 3GB RAM (available upon request); disk & CPU utilization typically very low, i.e. <1%, may run CPDN!
 - only bottleneck is mysql connections when doing hot backups or intensive queries (i.e. rerun credit calcs), but MySQL is getting better every version
- Scheduler/Web Server – Dell PowerEdge 6850, two Xeon 2.4GHz CPU, 1GB RAM, also usually <<1%
- Upload Servers – federated worldwide, donated, so vary from “off the shelf” PCs to shared space on a large Linux cluster.

BOINC Caveats for Projects*

- ~2-4 full-time weeks “learning curve” of playing with demo apps, server configs etc before any meaningful use
- BOINC is still “dynamic & evolving” -- if a project doesn't have at least one person keeping tabs on changes & upgrades (& maybe helping dev BOINC?), they risk being left behind
- What may be a fix or “nice feature” for one project can break another project!
- “Simple upgrades” can end up hectic, i.e. New database fields, trying to merge new HTML pages etc.

**well it can't just all be gushing praise! :-)*

Grid-ifying BOINC for CPDN

- BOINC is grid-computing in the sense of distributed scientific computing (F&K: “The Grid”, 1999), although IBM WCG and Ian Foster now seems to “count it as grid”
- CPDN also adds the “normal” grid element of the distributed datasets that require access by researchers
- BOINC doesn't really address this, i.e. **SETI@Home** scientists are the only ones who need (or are given) access to the finished runs
- CPDN will provide data via grid-enabled web services to such providers as the NERC Data Grid <http://ndg.badc.rl.ac.uk/>
 - New hire at RAL and replacement at ComLab will enable NERC-DG interface



Educational Outreach

- CPDN has public education via the website, media, and schools as an important facet of the project
- Website has much information on climate change and related topics to the CPDN program.
- Schools are running CPDN and comparing results, especially during National Science Week (starts 12/3/04) with special events at U Reading
- Students will host a debate on climate change issues, compare and contrast their results etc.



- Currently focused on UK schools, but as projects added and staff resources are gained plan to expand to other European schools and US schools

Students at Gosford Hill School, Oxon viewing their CPDN model

Summary

- CPDN's move to BOINC has paid off in terms of...
 - ease of developing future apps for CPDN, especially as we have an ambitious schedule of 3 more apps (models) in the next 6-12 months
 - gaining users through “sharing” with other BOINC projects, also when CPDN down BOINC users can run other projects & vice-versa
 - allowing us to concentrate more on scientific/modelling issues (not so much user credits/stats/etc – all of which BOINC handles well)