

A presentation slide with a light blue gradient background. The title "How to Write Fast Code?" is centered at the top in a large, bold, dark blue font. Below the title, there are two green rounded rectangles: one labeled "Fast Platforms" and another labeled "Good Techniques", separated by a plus sign. Under "Fast Platforms", there is a bulleted list: "• Multicore platforms", "• Manycore platforms", and "• Cloud platforms". Under "Good Techniques", there is a bulleted list: "• Data structures", "• Algorithms", and "• Software Architecture". At the bottom left, the text "18-645 – How to Write Fast Code?" is shown, and at the bottom right, the number "2".

## Instructors

- **Jike Chong**
  - Email: jike.chong@sv.cmu.edu
  - SkypeID: chong.18645
  - Office Hours: Tuesdays (6:00 – 8:00pm CA) and by appointment
- **Ian Lane**
  - Email: lane@cs.cmu.edu
  - SkypeID: lane.18645
  - Office Hours: Tuesdays and Thursday (08:00 - 11:00am CA)

18-645 – How to Write Fast Code?

3

## Instructor: Jike Chong

- 10+ years working on **Fast Platforms**, 6 years researching **Good Techniques**
- Adjunct Professor**, ECE, Carnegie Mellon Silicon Valley
  - Directs the CUDA Teaching Center and the CUDA Research Center at CMUSV
- Head of Data Sciences**, SimplyHired.com
  - Using machine-learning algorithms in data analytics
- Prior work:
  - Three first-authored US Patents, one pending
  - Intel Research Labs, Sun Microsystems, Xilinx Research Labs, Electric Cloud, etc
  - Ph.D. from UC Berkeley, M.S. and B.S. for Carnegie Mellon University

18-645 – How to Write Fast Code?

4

## Instructor: Ian Lane

- Assistant Research Professor (CMU-SV, LTI, ECE)
  - 10 years working on *Speech Recognition, NLP, MT* and *Machine Learning*
  - Co-Directs the CUDA Teaching and Research Center at CMU-SV
  - Heavy user of distributed computing (Maui, Condor, SLURM, HADOOP):**
    - Large data, computationally expensive algorithms and limited computational resources
    - State-of-the-art Automatic Speech Recognition systems contain many components and can require over 40,000 CPU-hours to train (4.5 years!)
    - For effective throughput must also consider Network, Disk and Memory Access

**High-throughput Computing:**  
Create new technologies → Deliver practical systems

18-645 – How to Write Fast Code?

## Outline

- Why Fast Code Now?
- What does 100x speedup mean?
- Course structure and organization

18-645 – How to Write Fast Code?

## Why Fast Code Now?

**Need is driven by the applications...**  
...NOT by the availability of the platforms.

Recognition	Mining	Synthesis
What is a tumor?	Is there a tumor here?	What if the tumor progresses?

Dubey, "Recognition, mining and synthesis moves computers to the era of tera", Technology@Intel Magazine, February 2005

18-645 – How to Write Fast Code? 7

## RMS Appears Everywhere

- **Recognition:** Interpretations of the world (with models)
- **Mining:** Understanding of the world (discover hidden patterns with model)
- **Synthesis:** Anticipate outcomes in the world (using the models to predict)

Recognition	Mining	Synthesis
What is a <i>hedge</i> ? What is the <i>interest rate</i> ?	Is there a <i>hedging</i> opportunity here?	What if <i>interest rates</i> were to go up?

Chen et al, "Convergence of Recognition, Mining, and Synthesis Workloads and Its Implications", Proceedings of the IEEE, Vol. 96, No. 5, pp. 790-807, May 2008

18-645 – How to Write Fast Code? 8

## Norman's Gulf

- Donald Arthur Norman: a scholar in cognitive science, design and usability
  - Author of the book *The Design of Everyday Things*
- More computation is being dedicated to bridging Norman's Gulf
  - Making computers easier to use
- What are some examples you see?

A diagram illustrating Norman's Gulf. It features two green rounded rectangles side-by-side. The left rectangle contains the text "Computer's Model". The right rectangle contains the text "Human's Conceptual Model". Between them are three curved arrows forming a cycle: one arrow points from the Computer's Model to the Human's Conceptual Model, another from the Human's Conceptual Model back to the Computer's Model, and a third from the Computer's Model back to itself.

D. A. Norman, *The Design of Everyday Things*. New York: Currency-Doubleday, 1989.

18-645 – How to Write Fast Code? 9

## What is being valued today?

The diagram illustrates the shift in valuation across three categories:

- Hardware → Commoditized:** An image of a circuit board is shown next to the word "Hardware", followed by a blue arrow pointing to the word "Commoditized".
- Software → Open Sourced:** A screenshot of a terminal window displaying a "Hello World" C program is shown next to the word "Software", followed by a blue arrow pointing to the word "Open Sourced".
- Data → King!**: A screenshot of a Facebook profile page is shown next to the word "Data", followed by a blue arrow pointing to the word "King!".

18-645 – How to Write Fast Code? 10

## Outline

- Why Fast Code Now?
- What does 100x speedup mean?
- Course structure and organization

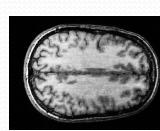
## What does 100x Speedup Mean?

**12 hours:**      10x speedup → 1.2 hours  
**100x speedup → 7.2 minutes**  
1000x speedup → 43.2 seconds

- Game changing technology advances
  - Overnight jobs becomes interactive



Speech Analytics

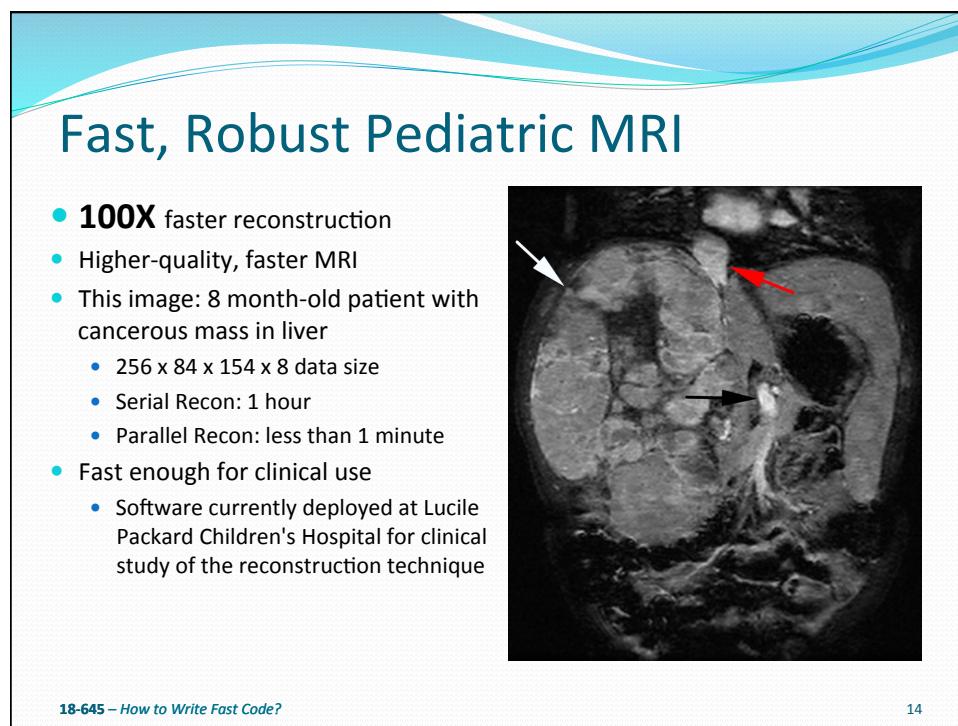
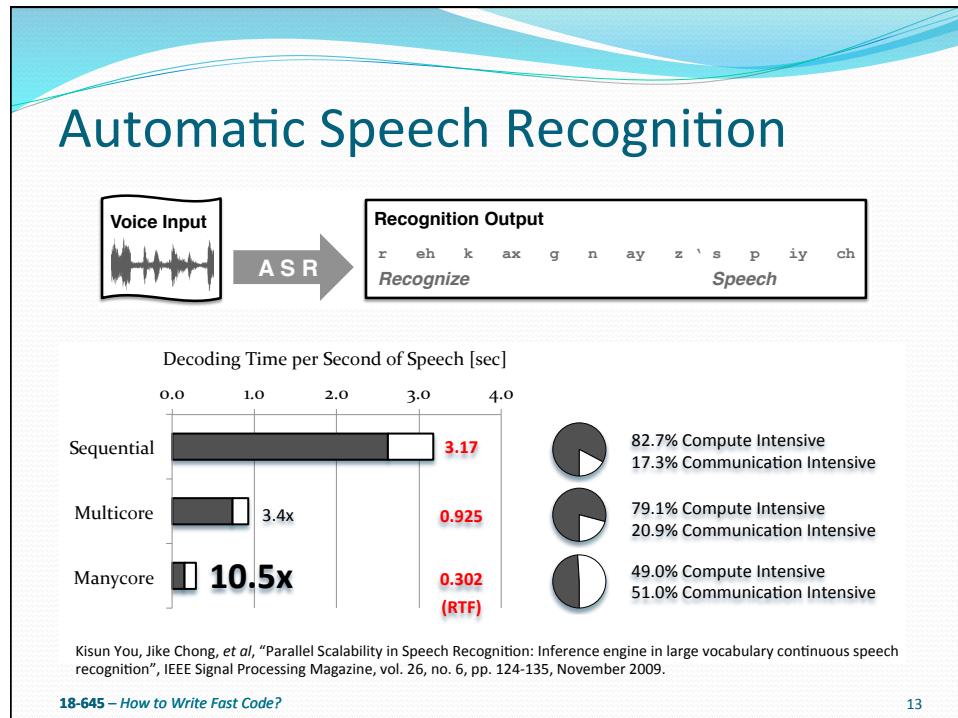


Medical Imaging



Image Recognition





## Support-Vector Machines

- Algorithmic changes and parallel implementation lead to performance speedup: Core 2 Duo versus G80

Computation	LIBSVM	Our algorithm 2-core Parallel CPU	Our algorithm, 16-core Parallel GPU
SVM Training (geo-mean)	771.8 s	---	38.5 s
SVM Classification	41.42 s	4.21 s	0.38 s

**100X speed-up**  
896 downloads since release in 10/2008

Fast support vector machine training and classification, Catanzaro, Sundaram, Keutzer, International Conference on Machine Learning 2008

18-645 – How to Write Fast Code?

15

## Image Contours Detection

Precision

Recall

CVPR 2008      Damascene

- We achieve equivalent accuracy on the Berkeley Segmentation Dataset
  - Comparing to human segmented “ground truth”
- F-measure 0.70 for both
- Human agreement = 0.79
- 3.8 minutes to 1.8 seconds: **126x speedup**
- 616 downloads since release in October 2009

"Efficient, High-Quality Image Contour Detection" Catanzaro, Su, Sundaram, Lee, Murphy, Keutzer International Conference on Computer Vision, 2009

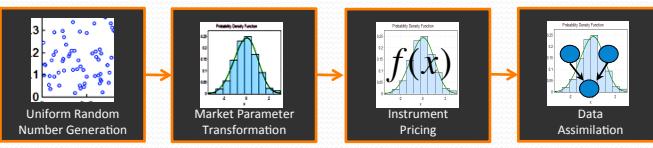
16

# Computational Finance

- Value-at-Risk Computation with Monte Carlo Method
- Summarizes a portfolio's vulnerabilities to market movements
- Important to algorithmic trading, derivative usage
- Improved implementation to run **60x faster** on a parallel microprocessor



Four Steps of Monte Carlo Method in Finance



Matthew Dixon, Jike Chong, Kurt Keutzer, "Acceleration of Market Value-at-Risk Estimation", Workshop on High Performance Computing in Finance at Super Computing 2009, November 15, 2009.

18-645 – How to Write Fast Code?

17

# NY Times TIFF to PDF

- In 2007, the New York Times decided to make all the public domain articles from 1851-1922 available free of charge
  - Needed to convert from TIFF to PDF

**Data set:** 4TB of raw image TIFF to 11 million PDF

**Compute Instances:** 100 Amazon EC2 instances

**Time taken:** 24 hours

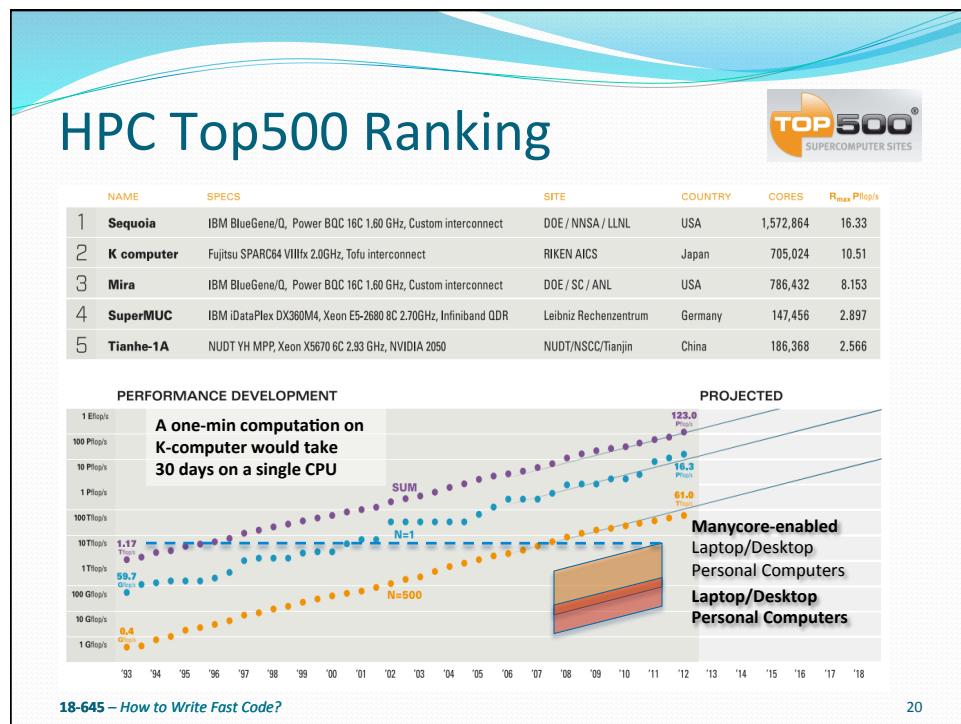
**Cost:** \$240



<http://open.blogs.nytimes.com/2007/11/01/self-service-prorated-super-computing-fun/>

18-645 – How to Write Fast Code?

18

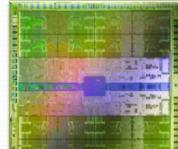


## Philosophy on Platforms

Provide theoretical background and hands-on practices...  
...to innovate with multicore/manycore/cloud based platforms.



Intel Sandy Bridge Multicore  
Processor (Core i7-2600K)



NVIDIA Fermi Manycore  
Processor - GTX580

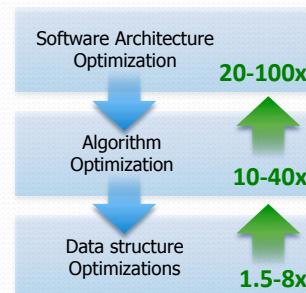


Yahoo! Hadoop Cluster  
~2000 nodes in one cluster

- **Significant scientific advances** will be empowered by multicore/manycore/cloud based platforms over the next decade
- Knowledge of these new computing capabilities will command a **dominating advantage** over your peers...  
...in developing innovative techniques to solve challenging problems

## Philosophy on Techniques

- Efficient **software architecture** is the most important to designing fast code
  - Software design patterns
- Understanding the implementation platform will help reason about application performance bottlenecks
- **Hands-on experience** provide confidence for you to effectively use these technologies for your research and development needs



## Outline

- Why Fast Code Now?
- What does 100x speedup mean?
- Course structure and organization

18-645 – How to Write Fast Code?

23

## Secret Word to Earn You a Spot!

- Secret word: **FLIRT**
- Send us an email with:
  - Subject: “**L01: Code FLIRT – Lastname, Firstname**”
  - To: “**18645@sv.cmu.edu**”
- Please use the same name as shown on your student record
- ***Send it now!***
- We will take the time stamp of this email into account when managing the waitlist

18-645 – How to Write Fast Code?

24

## Organization of This Course

- Bi-coastal broadcast
  - Every Tuesdays and Thursdays
  - Silicon Valley (Blg 23, Rm 118) 12:00 noon-1:20 pm PDT
  - Pittsburgh (HH1108) 3:00 pm-4:20 pm EDT
- Class website:
  - [fast645.info](http://fast645.info)
- During class time:
  - Lectures, Team presentations
  - Q&A sessions, Midterm/Final Exams
- Outside class time:
  - Three homework assignments
  - Three mini homework projects
  - One term project (focused on your own research or interests)

18-645 – How to Write Fast Code?

25

## Grading and Expectations

- GRADING
  - 10% Homework assignments (Assignments to help with the projects)
  - 10% Midterm examination
  - 30% Homework projects (Three equally weighted mini projects)
  - 20% Term Project
  - 30% Final examination
- EXPECTATIONS
  - Attend majority of lectures
  - Hand in all assignments and mini-projects
  - Complete a term project
  - Attend the final exam

18-645 – How to Write Fast Code?

26

## How to Earn Extra Credit!

- **Purpose:**
  - Help peers, do good work, get rewarded!
  - Reward can push up grade by  $\frac{1}{2}$  step: e.g. A to A+
- **Three Reward Opportunities:**
  - Use the Forum to get help and provide help from each other for setup the infrastructure for projects
    - Most useful answer to a discussion forum question gets **1** credit points
  - Mini homework projects will be tested **daily** for performance with **MOTION CHART!**
    - Top two performing team will get **2** credit points for each person for each day they are on top
  - Outstanding Final projects
    - Top 2-4 projects will be awarded **20** credit points for presenting at the HPC and GPU Supercomputing Group of Silicon Valley

18-645 – How to Write Fast Code?

27

## Software Reuse and Plagiarism

- **Question:** Can I learn to use code from the web?
- **Answer:** Yes! But only under certain circumstances...

Plagiarism	Acceptable SW Reuse
Collusion Unacknowledged code from a third-party	File-level Reuse Third-party code factored out to a separate file
Reverse Engineering Unacknowledged re-use of some code abstraction	Acknowledgement in Documentation Third-party code clearly distinguished in documentation
Translation Unacknowledged re-use by translating functions from one language to another	Adequate Testing Third-party code must be tested to one's own requirements
Code Generation Unacknowledged assistance by using some code generators	
Reuse Without Testing Reuse third-party code without testing against one's own requirements	

Gibson, "Software Reuse and Plagiarism: A Code of Practice", In Proceedings of ITiCSE'2009. pp.55-59

18-645 – How to Write Fast Code?

28

## Structure of Lectures

- **Module 1: Background**
  - Hardware, applications
- **Module 2: Multicore Programming**
  - Focus on application design with OpenMP
- **Module 3: Manycore Programming**
  - Focus on application design with CUDA
- **Module 4: Cluster Programming**
  - Focus on application design with Hadoop
- **Module 5: Special Topics**

## In Each Lecture

- We will provide a set of questions
  - You should be able to answer them after the lecture
- Today's question:
  - What does it mean to **FLIRT** with parallel computing technologies and working with the principles and practices of Writing Fast Code?
- You will be asked to answer those questions at the end of the class

## Structure of Team Presentations

- Two teams with the best performance will present their code
- Two teams with the slowest performance will have their code presented
- Multi-locational teams gets 5% extra points for project
- Best presentations will be invited to present at the HPC and GPU Supercomputing Group of Silicon Valley
  - Members representing engineers and managers at major Silicon Valley companies such as: Google, Apple, Intel, NVIDIA, AMD
  - A great opportunity to make connections for cool internships

## Structure of Q&A Sessions and Exams

- In-class Q&A Sessions
- Midterm:
  - Covers background, multicore/manycore application development
- Final:
  - Covers all material in the course

## Structure of Projects

- Three mini projects throughout the semester
  1. Multicore project with OpenMP
  2. Manycore project with CUDA
  3. Cloud project with Hadoop
- Projects are done in teams of two or three students
  - Where possible try and form teams across geographic locations
- Projects use **Git** code repositories
  - Project files must be checked out, optimized, and checked in
- Term project
  - An application area of your choice
  - One of the suggested projects
  - Term project introduction on

18-645 – How to Write Fast Code?

33

## Structure of Projects

- Programming assignments starting with fully working kernels
- Accelerate them with techniques discussed in class
- Automated scripts will pick up each teams project everyday
  - cron job starts at: 8pm PDT, 11pm EDT
  - Top three teams each day will accumulate bonus points
    - Can be used to push up a grade-level in borderline cases
- Write-up and presentations during class

18-645 – How to Write Fast Code?

34

## Structure of Homework

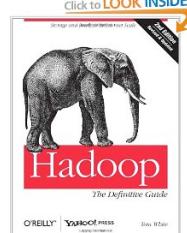
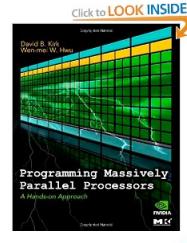
- Purpose:
  - Internalize concepts discussed in class
  - Opportunity to setting up the project environment
- Two week's time to complete the assignment
- Discussion among peers is allowed and encouraged
- However, tasks assigned and write-ups must be completed individually

18-645 – How to Write Fast Code?

35

## Recommended Books

- Required:
  - Kirk, Hwu, Programming Massively Parallel Processors, Second Edition: A Hands-on Approach, Morgan Kaufmann, 2012
- Optional:
  - White, Hadoop: The Definitive Guide, 3rd Edition, Yahoo Press, 2012



18-645 – How to Write Fast Code?

36

