

Deep Learning for Big Code

Martin Vechev
Spring 2018



Today

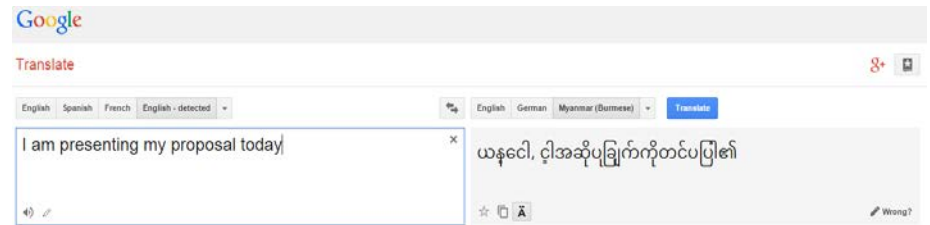
A bird's view of the area

Course organization

objectives, your tasks, grading, howto's

Learning and probabilistic models based on Big Data have revolutionized entire fields

Natural Language Processing
(e.g., machine translation)

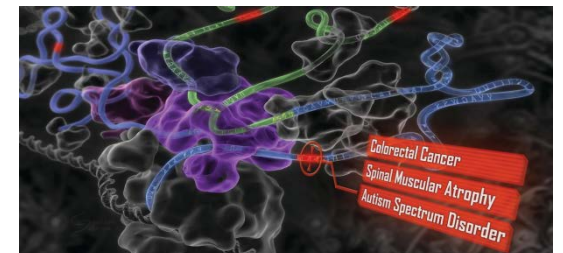


Computer Vision
(e.g., image captioning)



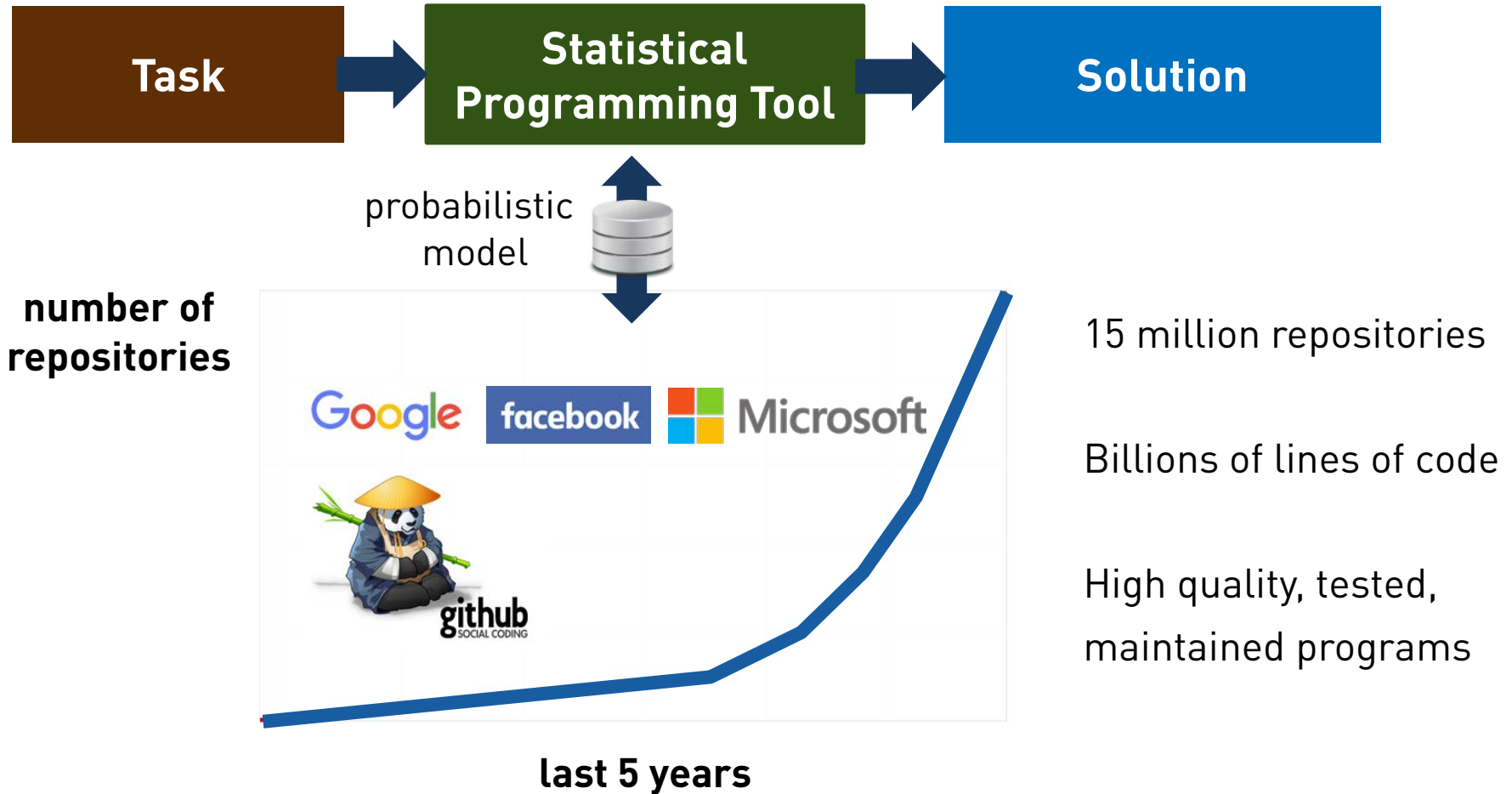
A group of people shopping at an outdoor market. There are many vegetables at the fruit stand.

Medical Computing
(e.g., disease prediction)



Can we bring this revolution to programmers?

Machine Learning for Programs



Why now?

Advances in Programming Languages
[Automated Reasoning, Synthesis, Constraint Solving]

Confluence of streams

Advances in Machine Learning
[Deep Learning, Graphical Models, Language Models]

Data
[> 15 million public repositories]



**machine learning-based
programming tools**
new rules, new ideas, new opportunities

Sample efforts in this space

Academia

@ETH: <http://plml.ethz.ch>

@Technion: <http://www.cs.technion.ac.il/~yahave/prime/>

@Rice: <https://www.cs.rice.edu/~sc40/>

@MIT: <http://people.csail.mit.edu/fanl/>

@Microsoft:
<https://www.microsoft.com/en-us/research/project/program/>

@Edinburgh: <http://homepages.inf.ed.ac.uk/csutton/>

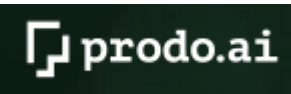
Start-ups

 DEEP CODE

<http://deepcode.ai> [ETH spin-off]

 codota

<http://codota.com>

 prodo.ai

[https://prodo.ai/](https://prodo.ai)

 NEAR

<http://near.ai>

DARPA's funded the 40M USD MUSE program on this topic

Seminal work in the area

**Learning from Large Codebases,
PhD Thesis, ETH Zurich, 2016**

ETH Medal for Outstanding PhD thesis

ACM Doctoral Dissertation, Honorable Mention Award

Only 3rd time in 40-year history of ACM that a PhD from Europe wins this award.

“...regarded as having the potential to open up several promising new avenues of research in the years to come...”

source: <https://awards.acm.org/about/2016-doctoral-dissertation>



Dr. Veselin Raychev
(CTO, DeepCode)

ACM SIGPLAN and Communications of the ACM Research Highlights, 2018

“...significantly advance the state-of-the-art in statistical reasoning on programs and offer the first concrete evidence of the tremendous promise of the overall approach...”

source: <http://www.sigplan.org/Highlights/Papers/>

Probabilistically likely solutions to problems hard to solve otherwise

Publications

- Program Synthesis for Character Level Language Modeling, **ICLR'17**
- Learning a Static Analyzer from Data, **CAV'17**
- Statistical Deobfuscation of Android Applications, **ACM CCS'16**
- Probabilistic Mode for Code with Decision Trees, **ACM OOPSLA'16**
- PHOG: Probabilistic Mode for Code, **ACM ICML'16**
- Learning Programs from Noisy Data, **ACM POPL'16**
- Predicting Program Properties from “Big Code”, **ACM POPL'15**
- Code Completion with Statistical Language Models, **ACM PLDI'14**
- Machine Translation for Programming Languages, **ACM Onward'14**

ML Engines

apk-deguard.com



DEEP3

jsnice.org



SLANG

nice2predict.org



more: <http://plml.ethz.ch>

<http://deepcode.ai>

Dimensions:

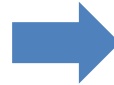
Machine Learning for Programming

Applications	Code completion Deobfuscation	Program synthesis	Feedback generation	Translation
Intermediate Representation	Sequences (sentences)	Trees	Translation Table	Graphical Models (CRFs) Feature Vectors
Analyze Program (PL)	typestate analysis scope analysis	control-flow analysis	alias analysis	
Train Model (ML)	Neural Networks N-gram/back-off	SVM PCFG	Structured SVM	
Query Model (ML)	$\operatorname{argmax}_{y \in \Omega} P(y x)$		Greedy MAP inference	

ML to write code

```
Camera camera = Camera.open();  
camera.setDisplayOrientation(90);
```

SLANG



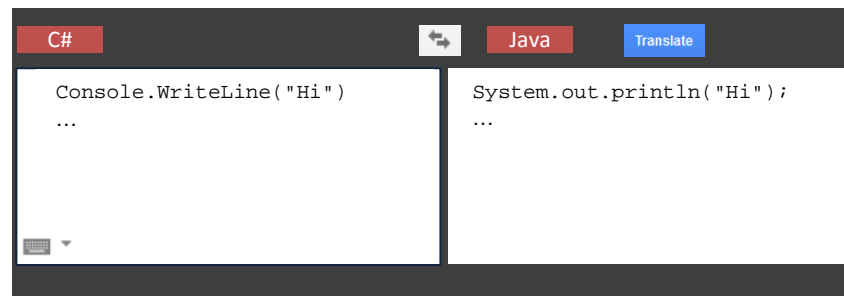
```
Camera camera = Camera.open();  
camera.setDisplayOrientation(90);  
  
camera.unlock();  
SurfaceHolder holder = getHolder();  
holder.addCallback(this);  
holder.setType(SurfaceHolder.STP);  
MediaRecorder r = new MediaRecorder();  
r.setCamera(camera);  
r.setAudioSource(MediaRecorder.AS);  
r.setVideoSource(MediaRecorder.VS);  
r.setOutputFormat(MediaRecorder.MPEG4);
```

Statistical language models
Recurrent neural networks

+

Typestate analysis
Alias analysis

ML to Translate Between Languages



The screenshot shows a web-based code translation tool. At the top, there are three tabs: 'C#' (selected), 'Java', and 'Translate'. Below the tabs are two text areas. The left text area contains C# code: `Console.WriteLine("Hi");` followed by an ellipsis. The right text area contains the translated Java code: `System.out.println("Hi");` followed by an ellipsis. A small keyboard icon is visible in the bottom-left corner of the C# text area.

Phrase-based Statistical
Machine Translation

+

Prefix Parsing of
Context-Free Grammars

ML for code deobfuscation

```
function FZ(e, t) { var n = [];  
var r = e.length; var i = 0;  
for (; i < r; i += t) if (i + t <  
r) n.push(e.substring(i, i +  
t)); else  
n.push(e.substring(i, r));  
return n;  
}
```



```
function chunkData(str, step) {  
  var colNames = [];  
  var len = str.length;  
  var i = 0;  
  for (; i < len; i += step)  
    if (i + step < len)  
      colNames.push(str.substring(i, i + step));  
    else  
      colNames.push(str.substring(i, len));  
  return colNames;  
}
```

JSNice.org

JS NICE

✓ Every country

✓ 200,000 users

✓ Top ranked tool



This Page Amsterdam @thispage_ams · Jul 16
Do you write ugly JavaScript code? Not to worry. JSNice will make it look like you are a **superstar** coder. Yai! - buff.ly/1HR4JL7

Ingvar Stepanyan @RRverser · Aug 6
JSNice.org became my **must-have** tool for code deobfuscation.
Expand

Brevity @seekbrevity · Jul 28
JSNice is an **amazing** tool for de-minifying #javascript files. JSNice.org, its great for #learning and reverse engineering.
Expand

Alvaro Sanchez @alvasavi · Jun 10
This is gold.
Statistical renaming, Type inference and Deobfuscation.
jsnice.org
Expand

Alex Vanston @mrvdot · Jun 7
I've been looking for this for years: JS NICE buff.ly/1pQ5qfr #javascript #unminify #deobfuscate #makeItReadable
Expand

Kamil Tomšik @cztomsik · Jun 6
tell me **how this** works!
de-minify #jquery #javascript incl. args, vars & #jsdoc impressive! jsnice.org



COMMENTS

MOST RECENT



Frid Kometz on:

20 Online Code Editors and Tools for Developers

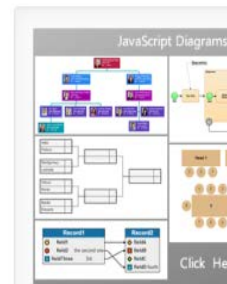


HOME WORDPRESS WEB DESIGN TYPOGRAPHY RESOURCES ARCHIVES CONTACT ADVERTISE

20 Essential Tools for Coders

BY GAVIN IN DEVELOPMENT RESOURCES -- 28 JUL, 2014

20 Best Freebies for Web Designers of Year 2014



After spending a lot of time in a particular field, we normally get a lot of experience and we suddenly start doing things easily and in short span of time to maximize our efforts towards our goal. This same procedure also holds true in case of web designing also. Because it's all about getting an experience in a particular field or language. Learning web designing is not an easy task, but when you hold experience in the same field, it becomes much easier for you.

AUGUST 28, 2014

JavaScript迷宮教星 — JS NICE

Even most of the popular JavaScript minifiers are not easy to use.

我們網頁工程師常常要在網站上用腳本debug開發網站上的JavaScript檔案時，通常那些檔案都已經經過最小化、最佳化、還有被加密了，以確保檔案運行的效率，並防止他人輕易的盜用程式碼。不過一旦要除錯，就苦了和編工程師了。

其實用「JS Nice」這個JSNice就是透過美化的程式碼功能，讓最小化的程式碼變成有條理的狀況，但是遇到腳本因為編譯程式，面對一些變數名稱如：b, c, 或程式，還是會覺得好煩好累。

不過網路上有些神奇的工具可以幫助我們不僅是將編譯一團亂的變數名稱，甚至它會替我們把編譯已經混淆的變數名稱！

這就是一個編譯介紹JSNice網站。

這個網站介紹如何乾淨、就像下圖一樣，在邊框裡輸入你的JS碼，只要輸入NICY JAVASCRIPT，經過解碼與除錯腳本的效果就會出現在右邊了！

right place, I have gathered 20 Essential JavaScript development tasks. Following development tasks, you will find the list handy and accurate.

1. JS Nice



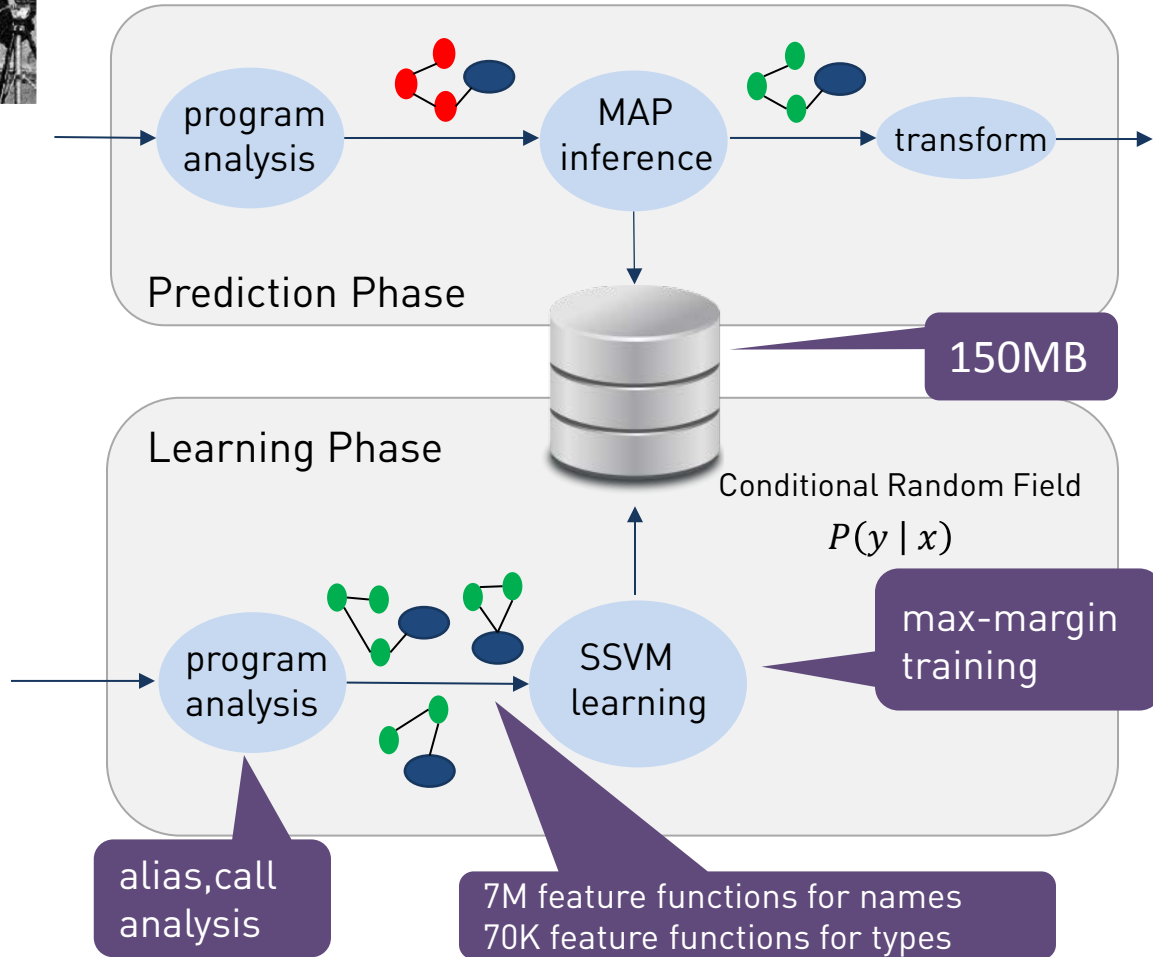
JS NICE STATISTICAL RENAMING, TYPE INFERENCE AND DEOBFUSCATION

ENTER JAVASCRIPT COPY JAVASCRIPT RESULT

1. // Put your JavaScript here that you want to rename, deobfuscate, ...

JSNice

```
var n = [];  
var r = e.length;  
var i = 0;  
for (; i < r; i += t)  
  if (i + t < r)  
    n.push(e.subs(i, i + t));  
  else  
    n.push(e.subs(i, r));  
return n;
```



```
var colNames = [];  
var len = str.length;  
var i = 0;  
for (; i < len; i += step)  
  colNames.push(str.subs(i, i + step));  
else  
  colNames.push(str.subs(i, len));  
return colNames;
```



ML for code deobfuscation

```
package a.b.c
class a extends SQLiteHelper
{
    SQLiteDatabase b; public
a(Context ctx) { b =
getWritableDatabase(); }
Cursor c(String str) {
return b.rawQuery(str); }}
```



```
package com.example.dbhelper

class DBHelper extends
SQLiteHelper {
    SQLiteDatabase db;

    public DBHelper(Context ctx) {
        db = getWritableDatabase();
    }

    Cursor execSQL(String str) {
        return db.rawQuery(str);
    }
}
```



Funny Reddit post/comment



[–] **Tycon712** • 3 points 2 days ago



Can someone tell me **what the point of using Proguard** is if there are tools out there like this?

[permalink](#) [embed](#) [pocket](#)



[–] **theheartbreakpug** • 6 points 2 days ago



As far as I know, this is brand new. I asked the creator of ProGuard a week ago how hard it is to unobfuscate code after it's run through proguard. He said it strips all the names out of the code so it's essentially impossible. **I'm super impressed by what they've done here.**

Nice2Predict.org:

scalable structured prediction framework

fully, **open sourced**,
Apache license

used by **various**
groups worldwide

JS NICE

DEGUARD

Your System
Here

NICE 2 Predict

Fast, Approximate
MAP inference

Fast, Parallel, Structured SVM
and Pseudo-Likelihood Training

Arbitrary factors and
indicator functions

OPPORTUNITY

Advances in Programming Languages
[Automated Reasoning, Synthesis, Constraint Solving]

Advances in Machine Learning

[Deep Learning, Graphical Models, Language Models]

Data 
[> 15 million public repositories]

Learning-based programming tools
new rules, new ideas, new opportunities

RICH PROBLEM SPACE

Applications	Code completion Deobfuscation	Program synthesis	Translation Feedback generation
Intermediate Representation	Sequences (sentences) Trees	Translation Table	Graphical Models (CRFs) Feature Vectors
Analyze Program (PL)	typestate analysis scope analysis	control-flow analysis	alias analysis
Train Model (ML)	Neural Networks N-gram language model	PCFGs	SVM Structured SVM
Query Model (ML)	$\text{argmax } P(y x)$ $y \in \Omega$		Greedy MAP inference

NEW PROBABILISTIC MODELS

1. **Pick** a structure of interest, e.g., trees:



2. **Define** a DSL for expressing functions:
(can be Turing complete)

```
TCond ::= ε | WriteOp TCond | MoveOp TCond
MoveOp ::= Up, Left, Right, DownFirst, DownLast,
NextDFS, PrevDFS, NextLeaf, PrevLeaf,
PrevNodeType, PrevNodeValue,
WriteOp ::= WriteValue, WriteType, WritePos
```

3. **Synthesize** $f_{best} \in \text{DSL}$ from Dataset D :

$$f_{best} = \underset{f \in \text{DSL}}{\text{argmin}} \text{cost}(D, f)$$

4. **Use** f_{best} on new structures:

$$f_{best} (\text{tree diagram}) \rightarrow \gamma$$

PRACTICAL IMPACT

Course objectives

Introduction to the emerging area of learning from Big Code.

Learn how to read and evaluate papers in the area

Learn how to make good technical presentations

Basic Information

- Instructors

- Prof. Dr. Martin Vechev (martin.vechev@inf.ethz.ch)
- Dr. Veselin Raychev (veselin.raychev@inf.ethz.ch)

- Meetings

- once a week, 2 presentations per meeting
- next meeting: March 5
- Presentation schedule will be posted on course

WWW: <http://www.srl.ethz.ch/bigcode18.php>

How it works: your tasks

select a paper,
get date

By this **Friday**, send e-mail to Veselin (veselin.raychev@inf.ethz.ch) with your **5 choices**

study paper

create presentation

meet advisor,
get feedback

give final presentation: 30 min
answer questions: 15 min

Participation:
ask good questions, attend all classes

Study paper

select a paper,
get date

study paper

create presentation

meet advisor,
get feedback

3 'C's of reading:

Carefully: lookup unknown terms,
read cited papers

Critically: find limitations, flaws

Creatively: think of improvements

Write: key ideas, **try examples** by hand

Consult with TA / Instructors, if questions
email also fine

give final presentation: 30 min
answer questions: 15 min

Participation:
ask good questions, attend all classes

Create presentation

select a paper,
get date

study paper

create presentation

meet advisor,
get feedback

Explain motivation for the work

Clearly present the technical solution and results

Use **your own** example, not the one in the paper

Outline limitations / improvements

Focus on the key/crucial concepts

Do not present all of the details

give final presentation: 30 min
answer questions: 15 min

Participation:
ask good questions, attend all classes

Meet advisor, get feedback

select a paper,
get date

study paper

create presentation

meet advisor,
get feedback

clear technical questions on the paper

get feedback on draft presentation

meeting mandatory a week or so before presentation

give final presentation: 30 min
answer questions: 15 min

Participation:
ask good questions, attend all classes

Grading

select a paper,
get date

study paper

create presentation

meet advisor,
get feedback

quality of your final presentation

how well you understood the material ?

how well you presented it ?

how well you answered the questions ?

we will take into account paper difficulty

participation

did you ask good questions ?

attendance will be taken (do not miss classes)

give final presentation: 30 min
answer questions: 15 min

Participation:
ask good questions, attend all classes

Useful Presentation Links

“How to give strong technical presentations”, by Markus Püschel

- <http://users.ece.cmu.edu/~pueschel/teaching/guides/guide-presentations.pdf>
- Please read the above slides !

“Even a geek can speak” , by Joey Asher

