

# **Technology Transfer of Machine Learning in Practice**

Ralf Herbrich  
Amazon

# Overview

- Background
- Technology Transfer Case Studies
  - TrueSkill: Gamer Rating and Matchmaking
  - Click-Through Rate Prediction in Online Advertising
- Technology Transfer Lessons
  - Process
  - Technical

# Overview

- **Background**
- Technology Transfer Case Studies
  - TrueSkill: Gamer Rating and Matchmaking
  - Click-Through Rate Prediction in Online Advertising
- Technology Transfer Lessons
  - Process
  - Technical

# Background



1992 – 1997 (Berlin, Diploma)



1997 – 2000 (Berlin, PhD)



2000 – 2009 (Microsoft Research)



2009 – 2011 (Microsoft)



2011 – 2012 (Facebook)



2012 – Present (Amazon)



# Technology Transfer

- **Definition (Wiki):** Process of moving promising research topics into a level of maturity ready for bulk manufacturing or production
- **Practice:** Often failing due to
  - Different Success Criteria (Product vs. Publication)
  - No Training Programs for Technology Transfer
  - Processes Are Hard to Generalize (Structure?)

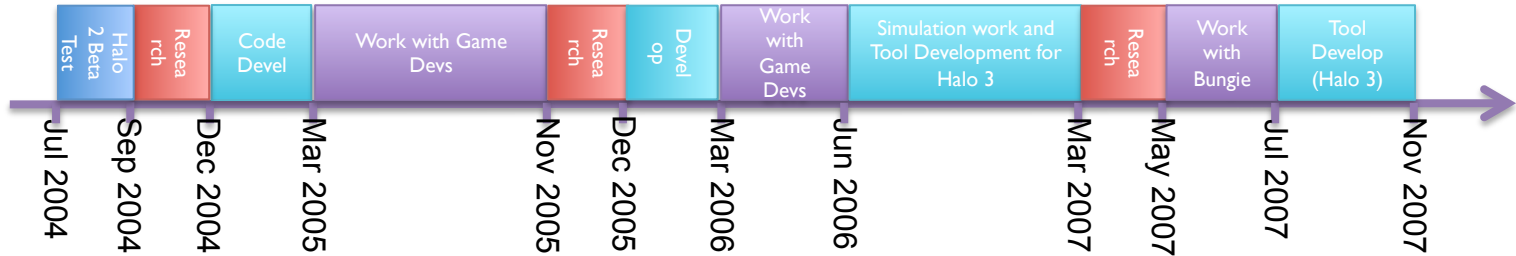
# Overview

- Background
- **Technology Transfer Case Studies**
  - TrueSkill: Gamer Rating and Matchmaking
  - Click-Through Rate Prediction in Online Advertising
- Technology Transfer Lessons
  - Process
  - Technical

# TrueSkill™

Joint work with Thore Graepel, Tom Minka & Phillip Trelford

# TrueSkill Technology Transfer



- **Lessons Learned:**

1. Pure research takes a short amount of time
2. Most of development was tool development
3. A platform feature only lives with a community
4. Mathematical Optimality  $\neq$  Fun Experience

# TrueSkill Technology Transfer



# The Skill Rating Problem

- **Given:**
  - Match outcomes: Orderings among  $k$  teams consisting of  $n_1$

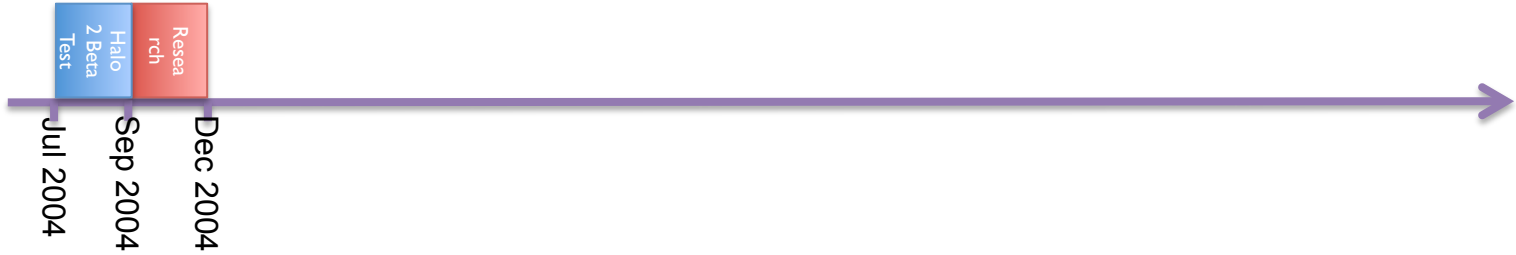
- **Que**
- **CL**

The image shows several overlapping screenshots from a game interface. At the top, a team score screen shows '1st Red Team' with a score of 50. Below it, a player stats screen for 'SniperEye' shows '1st' rank, 'N/A' level, and 'N/A' Avg. Life. A large red arrow points from this stats screen to a ranked list of players on the right. The list shows 17 players with their scores, ranging from 27 down to 23. The bottom of the image shows a partially visible screen with 'S3' and '7birds tm'.

Rank	Score	Player Name
1	27	SEWICSYDE OWNS
2	26	FATAL REVENGE
3	25	Paranoia 1
4	25	Paulk
5	25	IxX OMG Xxl
6	25	BittyTom
7	24	brian 2007
8	24	SEXY MOZES
9	24	droplates
10	24	jaCKdaSaMuRai
11	24	ll Me ll
12	24	iamNightMare
13	24	a retarded007
14	24	Perfected Brit
15	24	THE MUFFIN MANx
16	23	TheVunit
17	23	Mr Sushi87

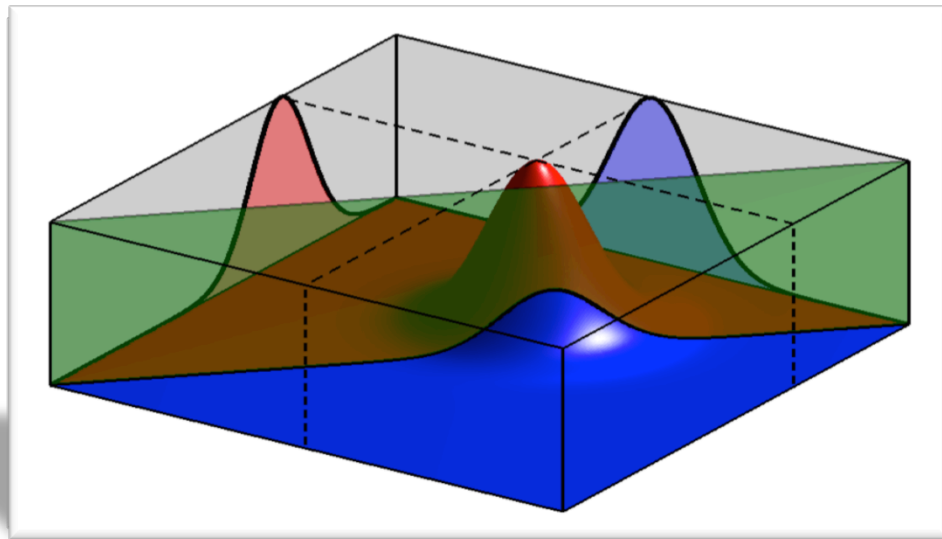
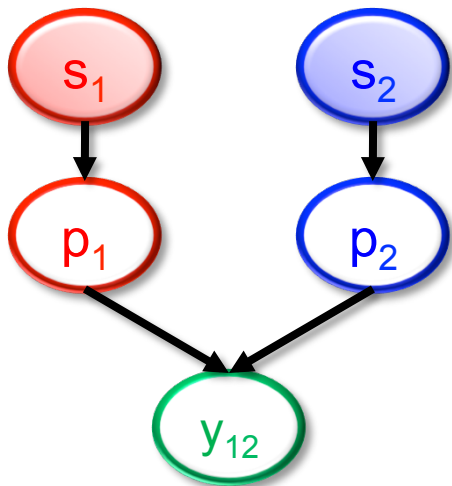
– All matches between teams of players

# TrueSkill Technology Transfer



# Two Player Match Outcome Model

- Latent Gaussian performance model for fixed skills
- Possible outcomes: Player 1 wins over 2 (and vice versa)

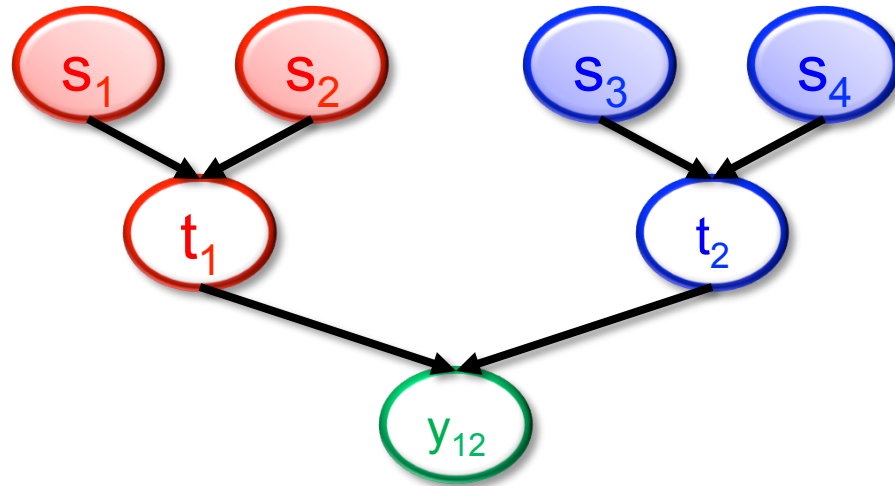


$$\mathbf{P}(y_{12} = (1, 2) | p_1, p_2) = \mathbb{I}(p_1 > p_2)$$



# Two Team Match Outcome Model

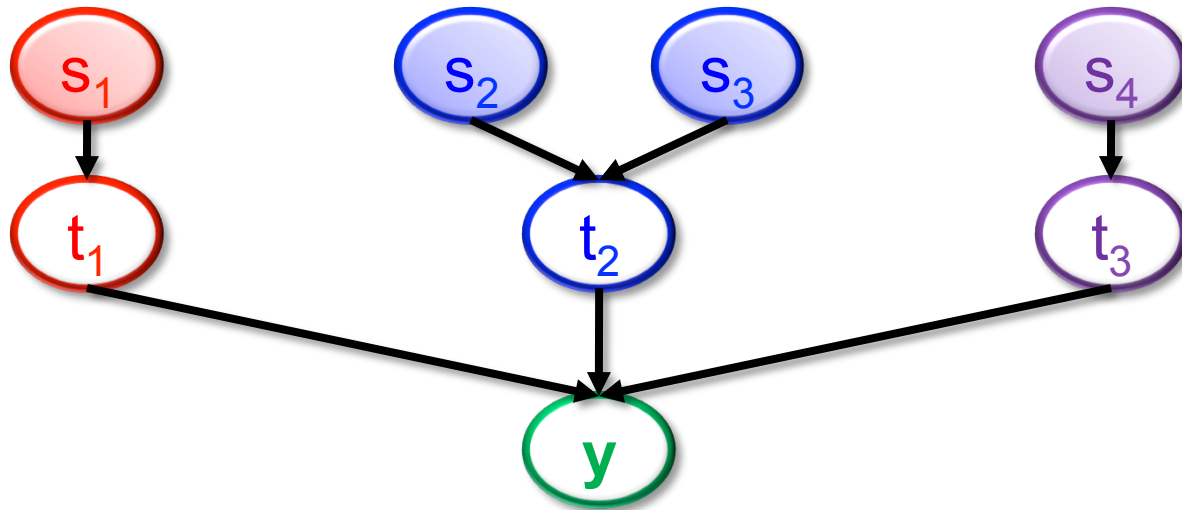
- Skill of a team is the sum of the skills of its members



$$\mathbf{P}(t_1 | s_1, s_2) = \mathcal{N}(t_1; s_1 + s_2, 2 \cdot \beta^2)$$

# Multiple Team Match Outcome Model

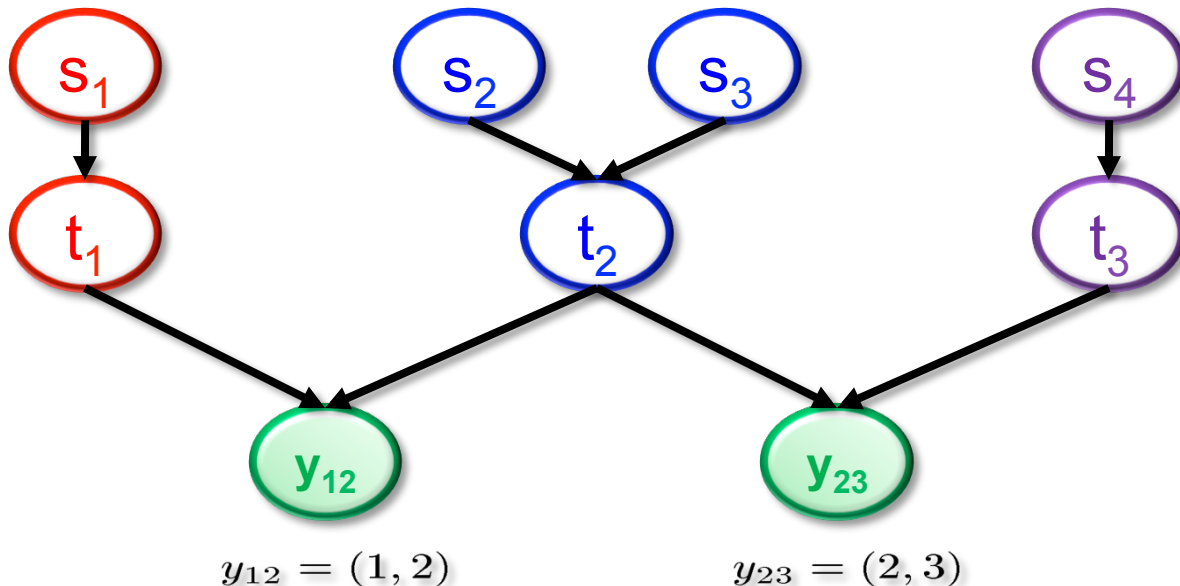
- Possible outcomes: Permutations of the teams



$$\mathbf{P}(\mathbf{y}|t_1, t_2, t_3) = \mathbb{I}(\mathbf{y} = (i, j, k)) \text{ where } t_i > t_j > t_k$$

# Multiple Team Match Outcome Model

- But we are interested in the (Gaussian) posterior!  $P(s_i | \mathbf{y} = (1, 2, 3)) = \mathcal{N}(s_i; \mu_i, \sigma_i^2)$



# Applications to Gaming

- **Leaderboard**

- Global ranking of all players

$$\mu_i - 3 \cdot \sigma_i$$

- **Matchmaking**

- For gamers: Most uncertain of
- For inference: Most informative

$$\frac{\mathbf{P}(p_i \approx p_j | \mu_i - \mu_j, \sigma_i^2 + \sigma_j^2)}{\mathbf{P}(p_i \approx p_j | \mu_i - \mu_j = 0, \sigma_i^2 + \sigma_j^2 = 0)}$$

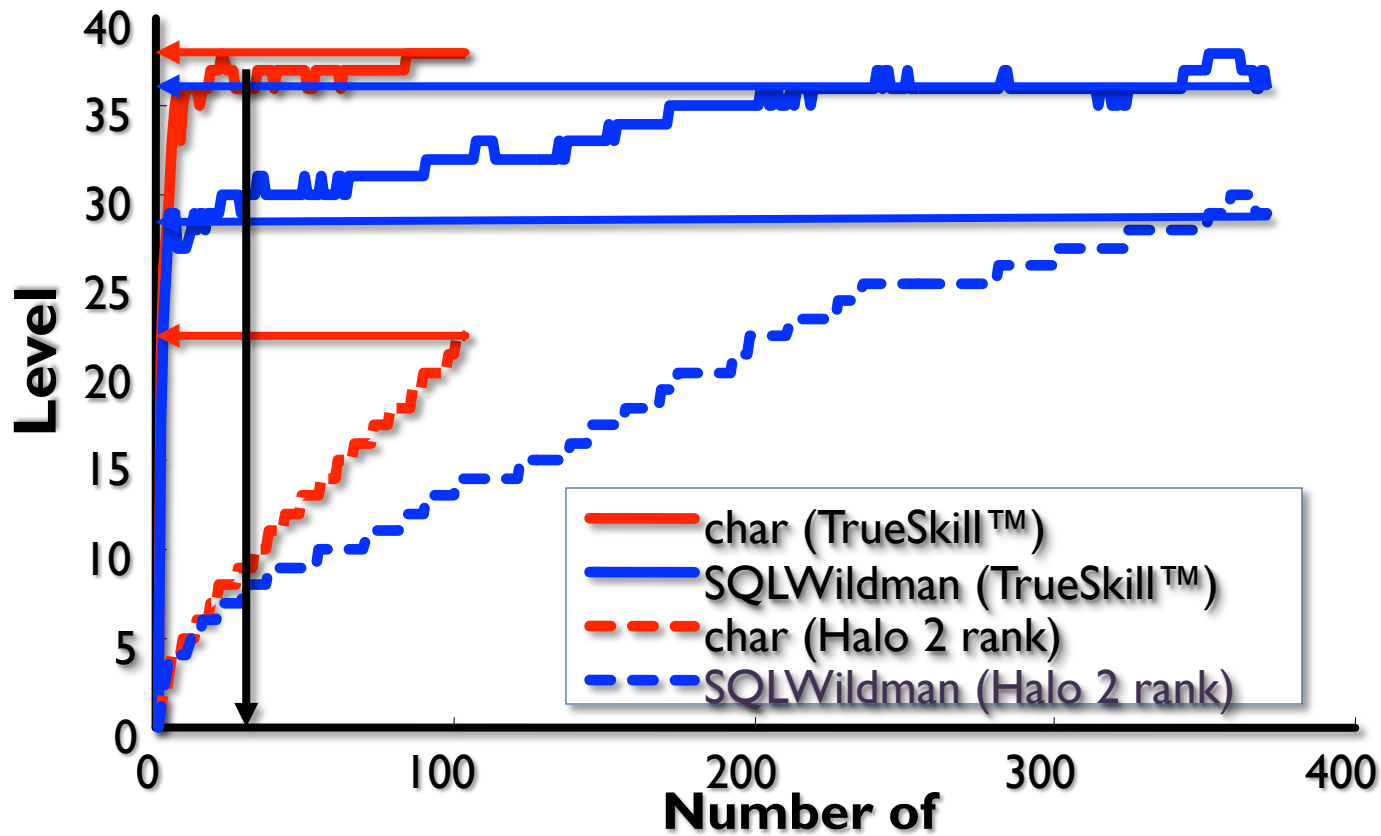
1	27	SEWiCSYDE OWNS
2	26	FATAL REVENGE
3	25	Paranoia 1
4	25	Paulk
5	25	IxX OMG Xxl
6	25	BittyTom
7	24	brian 2007
8	24	SEXY MOZES
9	24	droplates
10	24	jaCKdaSaMuRai
11	24	II Me II
12	24	iamNightMare
13	24	a retarded007
14	24	Perfected Brit
15	24	THE MUFFIN MANx
16	23	TheVunit
	23	Mr Sushi87

# Experimental Setup

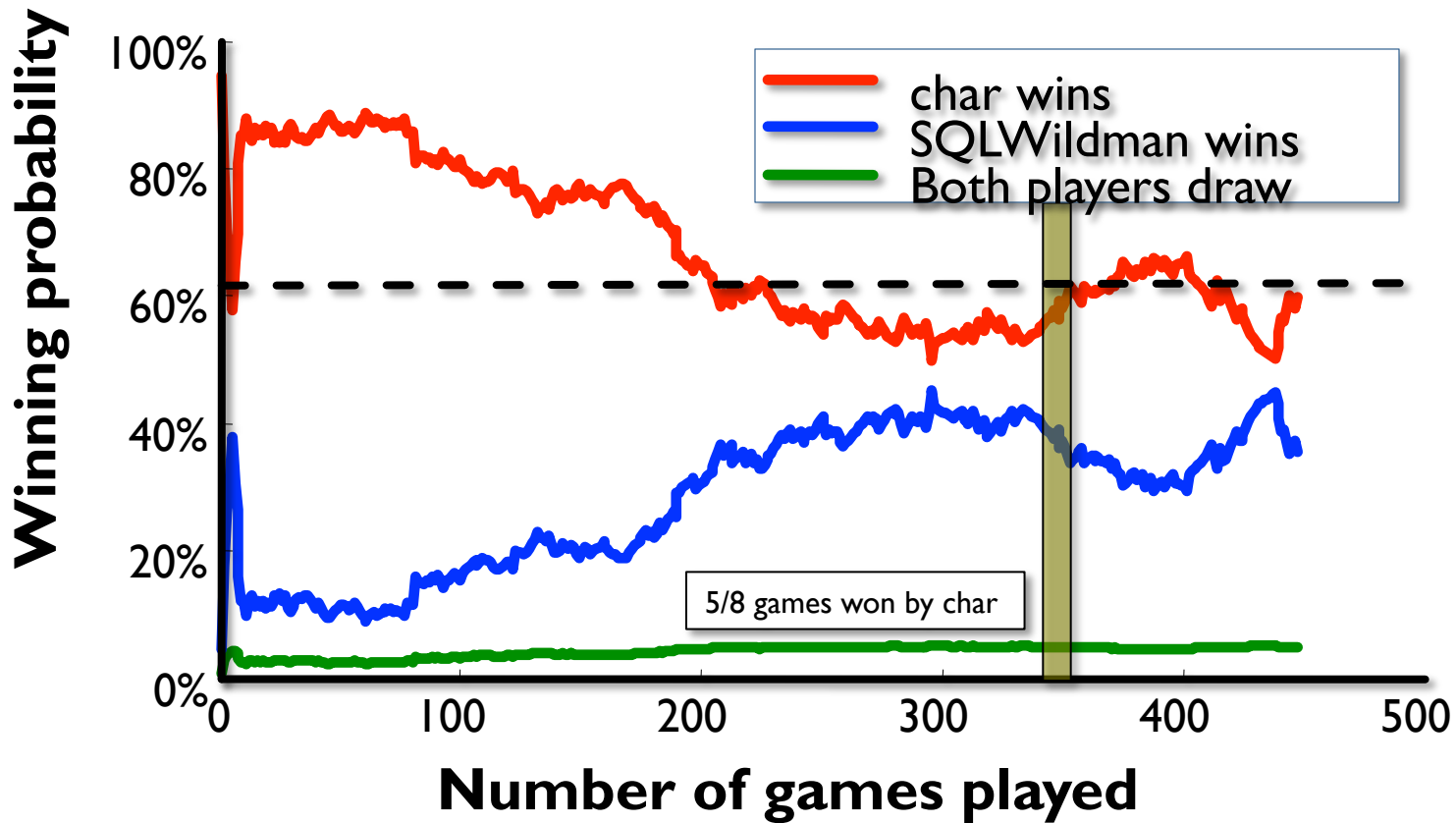
- **Data Set: Halo 2 Beta**
  - 3 game modes
    - Free-for-All
    - Two Teams
    - 1 vs. 1
  - > 60,000 match outcomes
  - $\approx$  6,000 players
  - 6 weeks of game play
  - Publically available



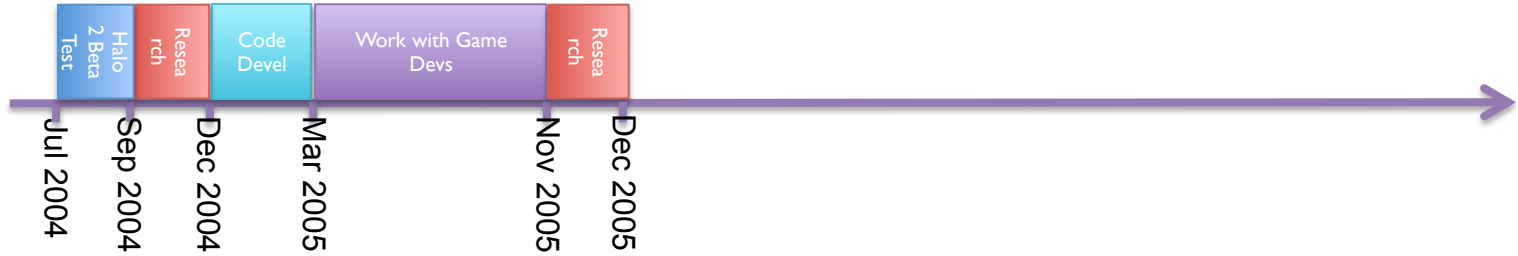
# Convergence Speed



# Convergence Speed (ctd.)




# TrueSkill Technology Transfer





# Graphical Models

- **Definition:** Graphical representation of joint probability distribution
  - Nodes:  = Variables
  - Edges: Relationship between variables
- **Variables:**
  - Observed Variables: Data
  - Unobserved Variables: ‘Causes’ + Temporary/Latent
- **Key Questions:**
  - (Conditional) *Dependency*:
  - *Inference*/Marginalisation:

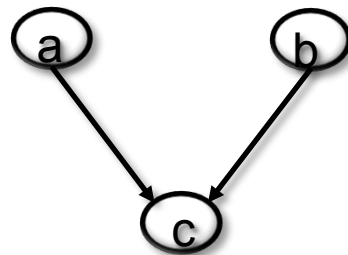
$$p(a, b|c) \stackrel{?}{=} p(a|c) \cdot p(b|c)$$
$$p(a, b) = \sum_c p(a, b, c)$$

# Directed Models: Bayesian Networks

- **Definition:** Graphical representation of joint probability distribution (Pearl, 1988)
  - Nodes: ○ = Variables
  - Directed Edges: Conditional probability distribution

- **Semantic:**

$$p(\mathbf{x}) = \prod_i p(x_i | \mathbf{x}_{\text{parents}(i)})$$



- Ancestral relationship of dependency

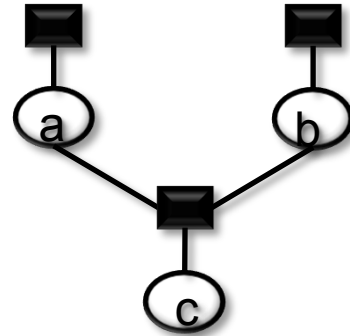
$$p(a, b, c) = p(a) \cdot p(b) \cdot p(c|a, b)$$

# Factor Graphs

- **Definition:** Graphical representation of product structure of a function (Wiberg, 1996)
  - Nodes:  $\blacksquare$  = Factors     $\bigcirc$  = Variables
  - Edges: Dependencies of factors on variables.

- **Semantic:**

$$p(\mathbf{x}) = \prod_f f(\mathbf{x}_{V(f)})$$



- Local variable dependency of factors

$$p(a, b, c) = f_1(a) \cdot f_2(b) \cdot f_3(a, b, c)$$

# Factor Graphs and Bayes' Law

- Bayes' law

$$p(\mathbf{s}|\mathbf{y}) \propto p(\mathbf{y}|\mathbf{s}) \cdot p(\mathbf{s})$$

- Factorising prior

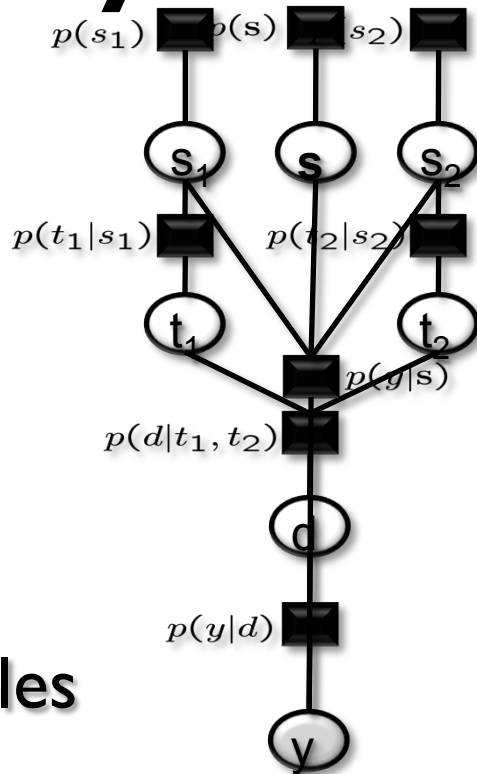
$$p(\mathbf{s}) = p(s_1) \cdot p(s_2)$$

- Factorising likelihood

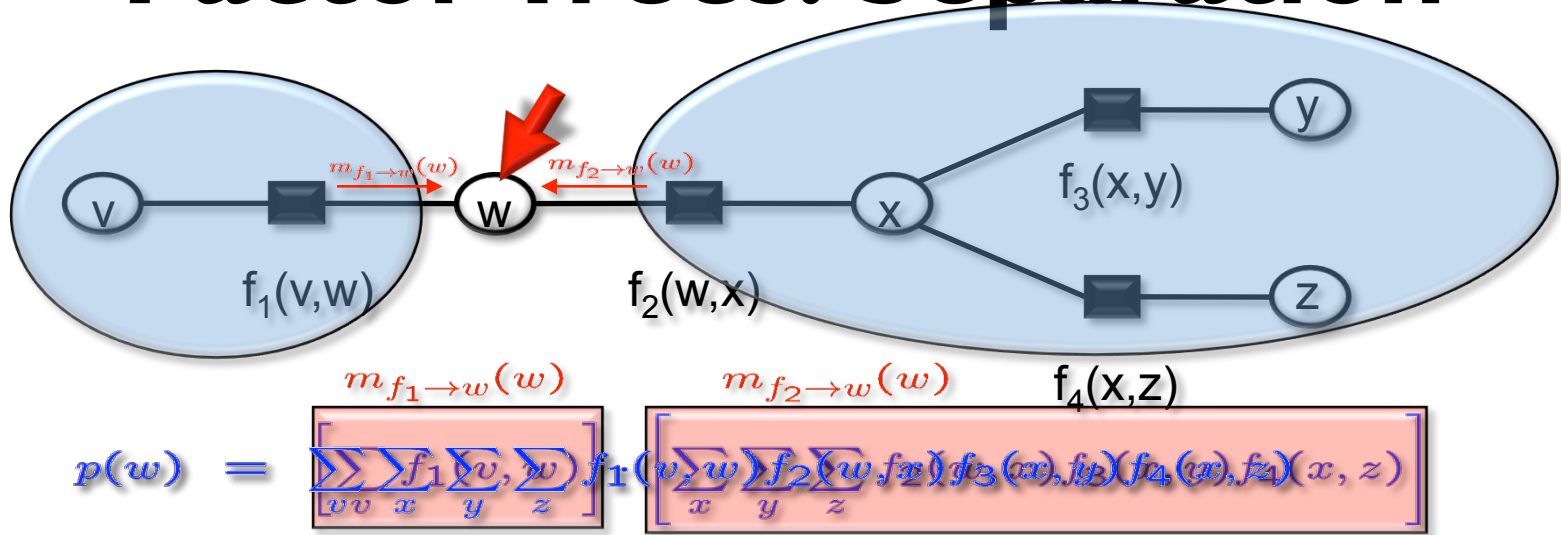
$$p(\mathbf{y}, \mathbf{t}, d|\mathbf{s}) = \prod_i p(t_i|s_i) \cdot p(d|t_1, t_2) \cdot p(y|d)$$

- Inference: Sum out latent variables

$$p(\mathbf{y}|\mathbf{s}) = \sum_{\mathbf{t}} \sum_d p(\mathbf{y}, \mathbf{t}, d|\mathbf{s})$$

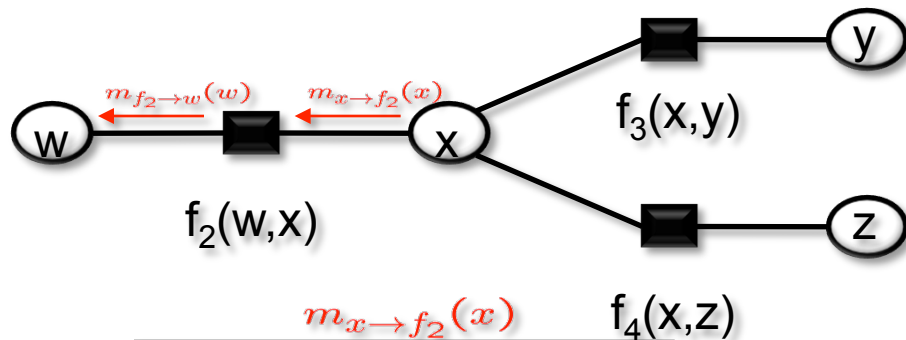


# Factor Trees: Separation



**Observation:** Sum of products becomes product of sums of all messages from neighbouring factors to variable!

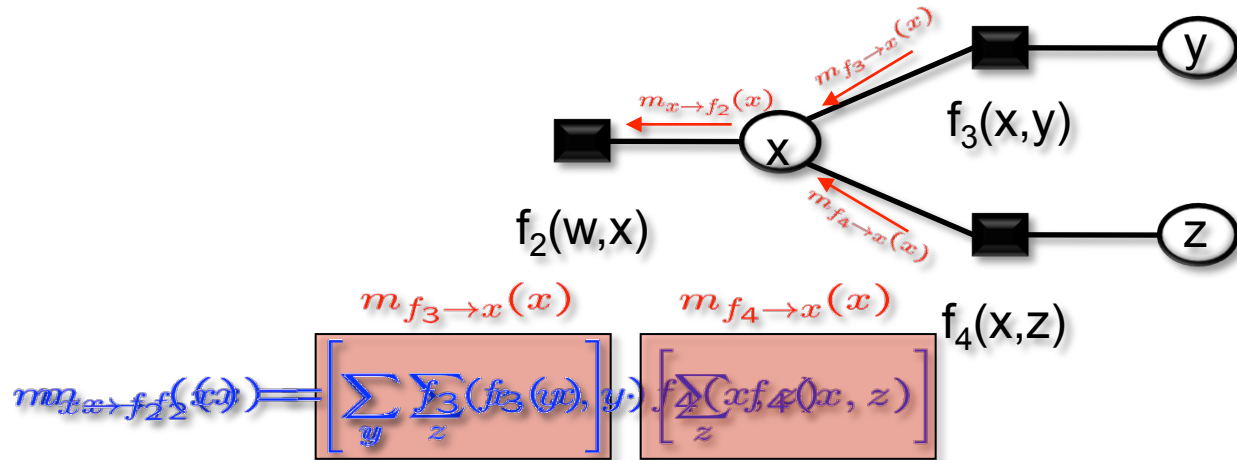
# Messages: From Factors To Variables



$$m_{f_2 \rightarrow w}(w) = \sum_x \sum_y \sum_z f_2(w, x) \left[ \sum_y \sum_z f_3(x, y) f_4(x, z) \right]$$

**Observation:** Factors only need to sum out all their local variables!

# Messages: From Variables To Factors



**Observation:** Variables pass on the product of all incoming messages!

# The Sum-Product Algorithm

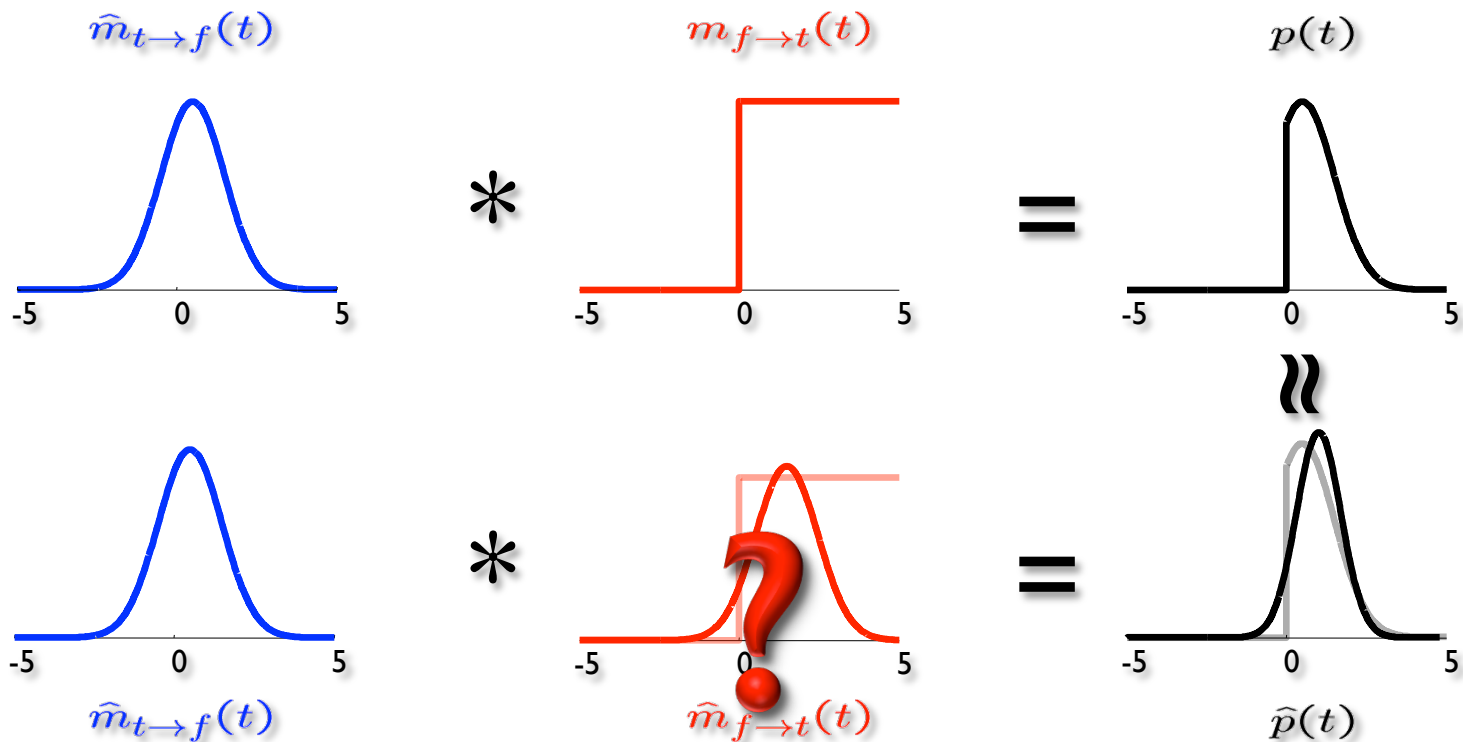
- Three update equations (Aji & McEliece, 1997)

$$p(t) = \prod_{f \in F_t} m_{f \rightarrow t}(t)$$
$$m_{f \rightarrow t_1}(t_1) = \sum_{t_2} \sum_{t_3} \cdots \sum_{t_n} f(t_1, t_2, t_3, \dots) \prod_{i > 1} m_{t_i \rightarrow f}(t_i)$$
$$m_{t \rightarrow f}(t) = \prod_{f_j \in F_t \setminus \{f\}} m_{f_j \rightarrow t}(t)$$

- Update equations can be directly derived from the distributive law.
- Calculate all marginals at the same time!
- Only need to pass messages twice along each edge!



# Approximate Message Passing



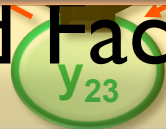
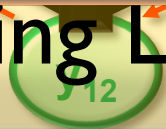
# Efficient Approximate Inference

## Gaussian Prior Factors

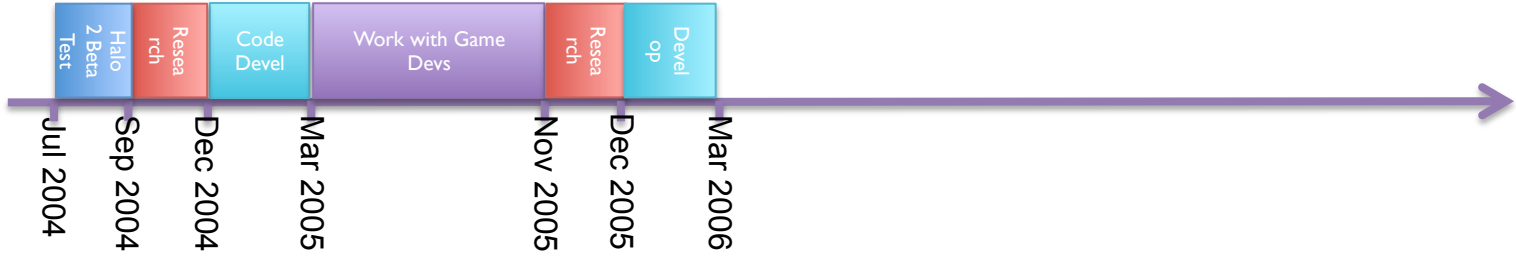
Fast and efficient approximate message passing  
using Expectation Propagation

[http://blogs.technet.com/b/apg/archive/  
2008/06/16/trueskill-in-f.aspx](http://blogs.technet.com/b/apg/archive/2008/06/16/trueskill-in-f.aspx)

Ranking Likelihood Factors

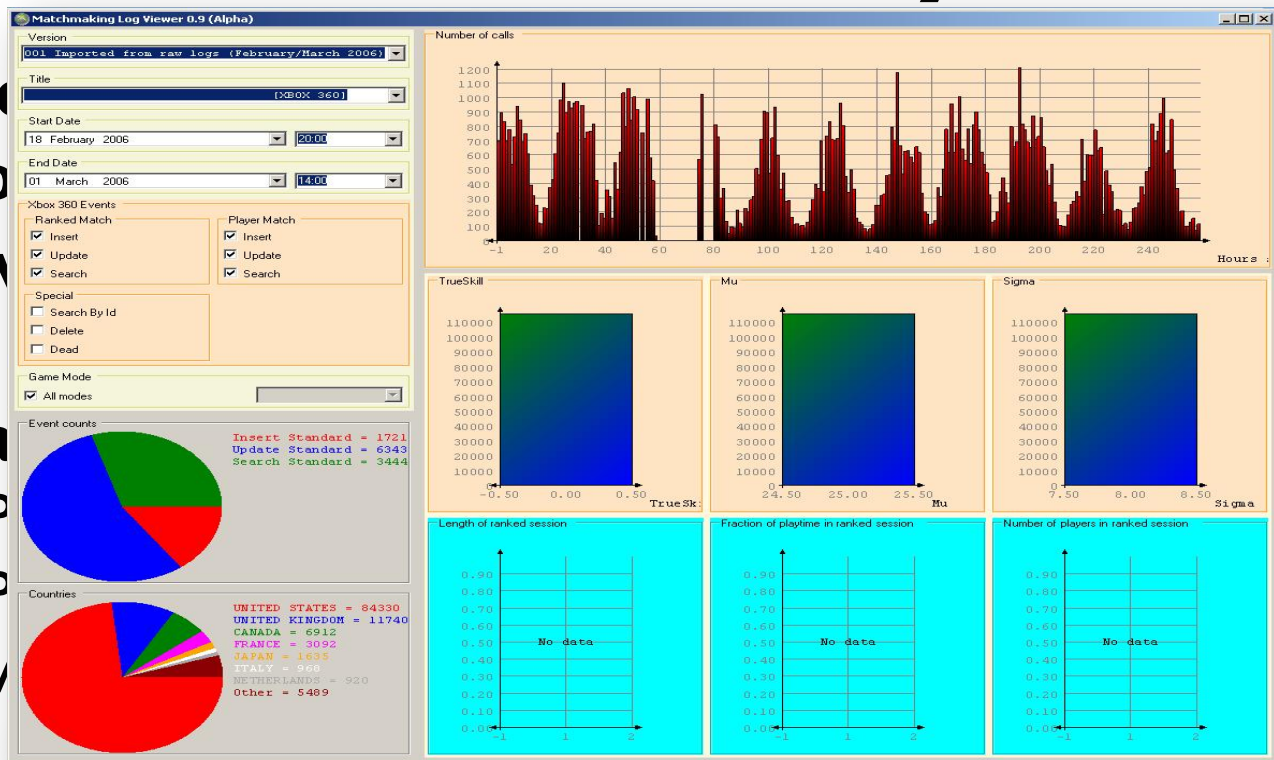


# TrueSkill Technology Transfer

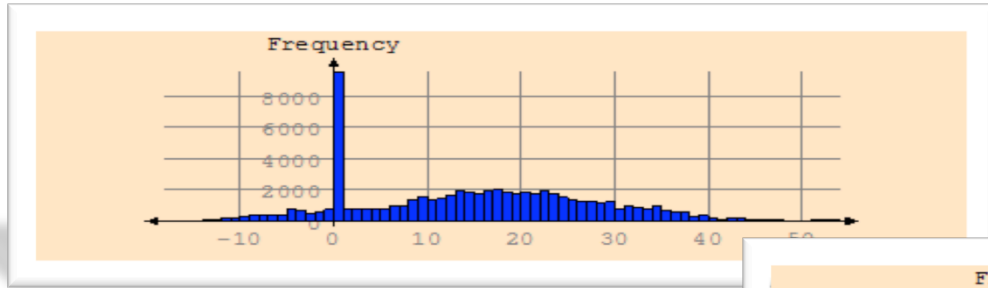


# Xbox Live Activity viewer

- Co
- Pro
- Dev
- Fea
- P
- P
- V

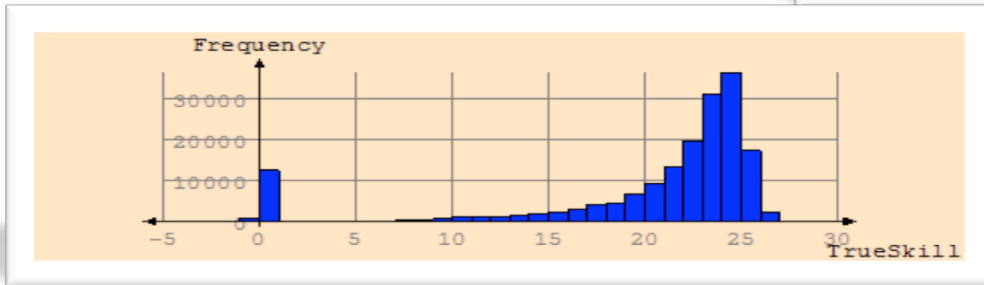
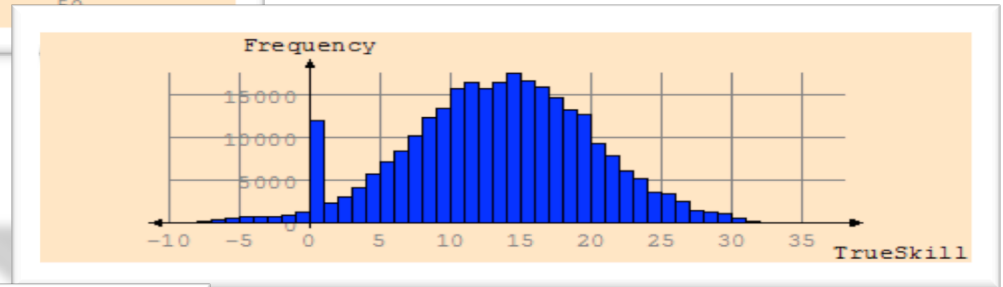


# Skill Distributions of Online Games



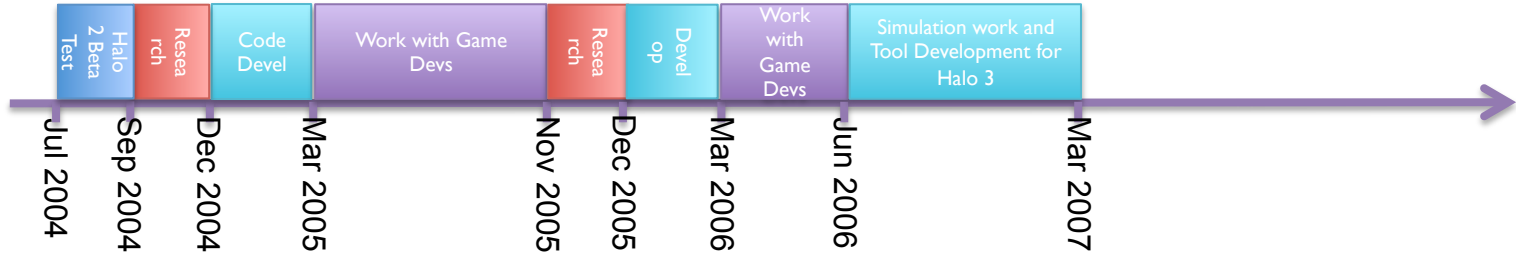
**Golf (18 holes): 60 levels**

**Car racing (3-4 laps): 40 levels**



**UNO (chance game): 10 levels**

# TrueSkill Technology Transfer



# Halo 3 in Action



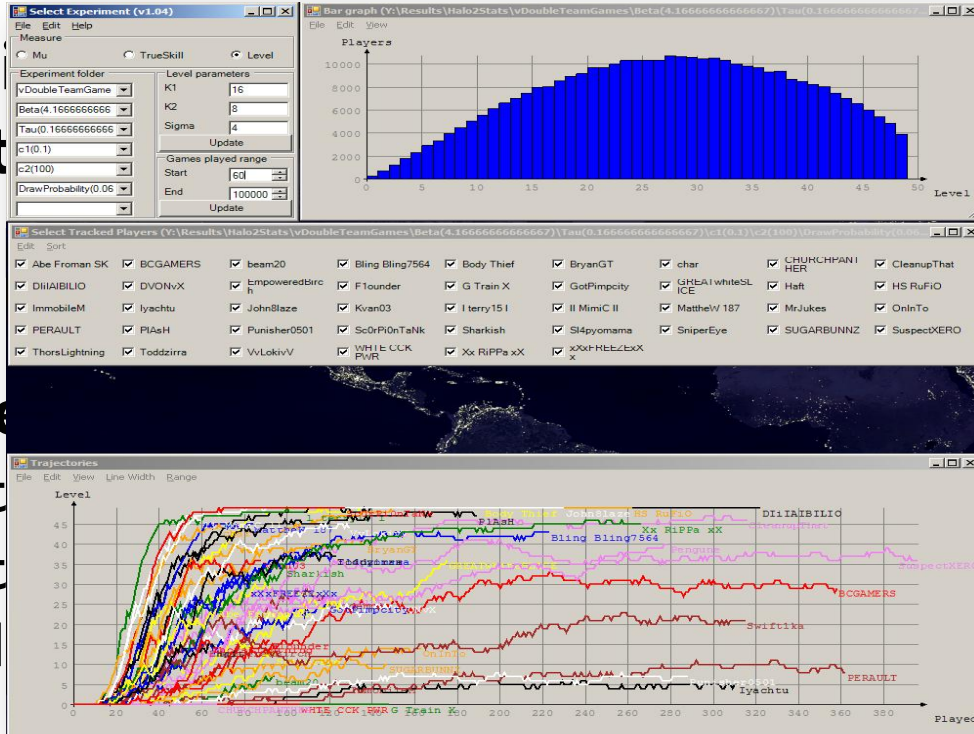
# Tools for Halo 3

- **Questions**
  - Controllable player skill progression (slow-down!)
  - Controllable skill distributions (re-ordering)
- **Simulations**
  - Large scale simulation of  $> 8,000,000,000$  matches
  - Distributed application written in C# using .Net remoting
- **Tools**
  - Result viewer (Logged results: 52 GB of data)
  - Real-time simulator of partial update



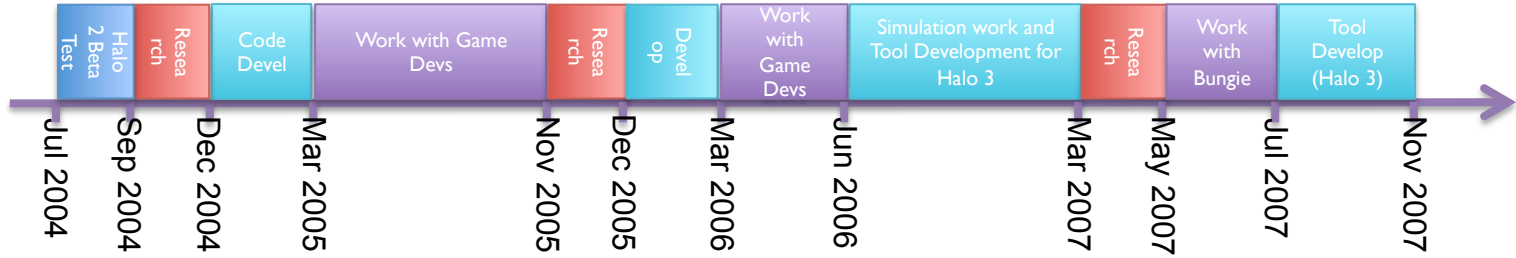
# Halo 3 Simulation Result Viewer

- Code s
- Project
- Develop
- Feature
  - Multi
  - Real-t
  - Based



(size)  
(easing)

# TrueSkill Technology Transfer



# Halo 3 Partial Update Analyser

- Code size: 2600 LOC

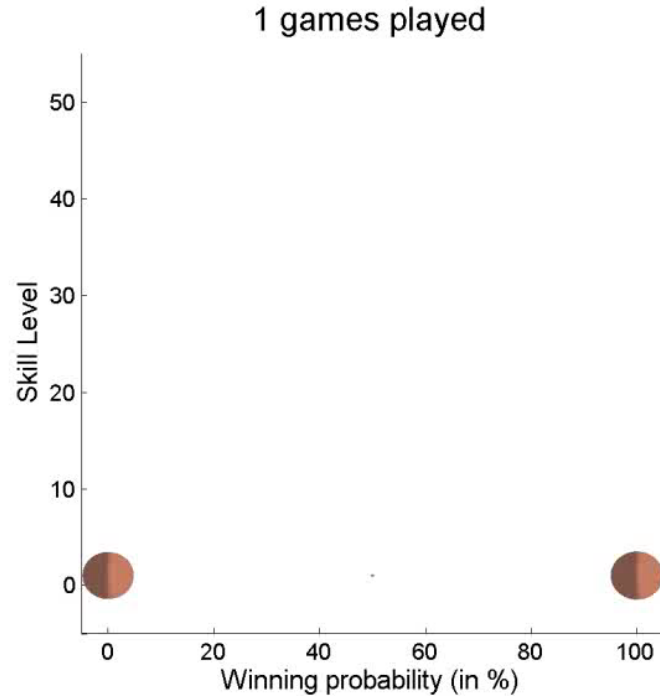
- P
- D
- F



Censored

- Real time changes

# Halo 3 Public Beta Analysis

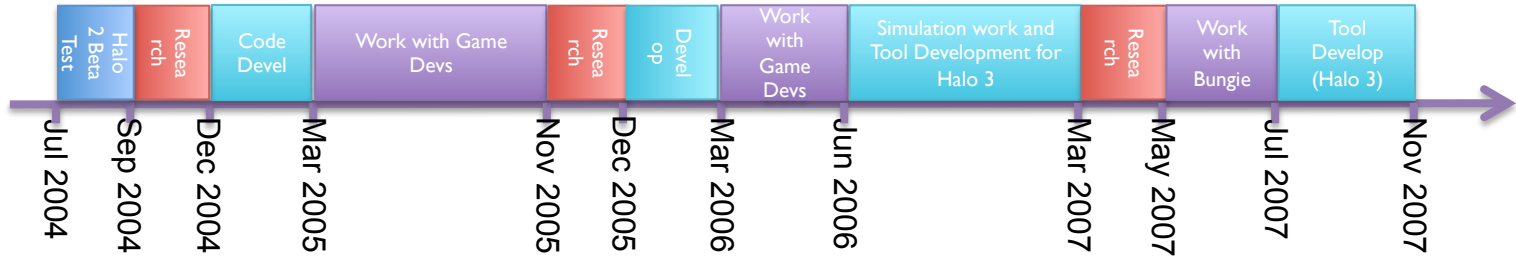


# Xbox 360 & Halo 3

- **Xbox 360 Live**
  - Launched in September 2005
  - Every game uses TrueSkill™ to match players
  - > 10 million players
  - > 2 million matches per day
  - > 2 billion hours of gameplay
- **Halo 3**
  - Launched on 25<sup>th</sup> September 2007
  - Largest entertainment launch in history
  - > 200,000 player concurrently (peak: 1,000,000)



# TrueSkill Technology Transfer



- **Lessons Learned:**

1. Pure research takes a short amount of time
2. Most of development was tool development
3. A platform feature only lives with a community
4. Mathematical Optimality  $\neq$  Fun Experience

# AdPredictor Technology Transfer



- **Lessons Learned:**

1. Pure research takes a short amount of time
2. Development takes much longer than planned
3. Counter-factual analysis and metrics are important
4. Develop for scale from Day 1

# adPredictor

Joint work with Thore Graepel, Joaquin Quiñonero Candela, Onno Zoeter, Tom Borchert , Phillip Trelford



# AdPredictor Technology Transfer



# Why Predict Probability-of-Click?

The screenshot shows a Windows Internet Explorer browser window displaying a Live Search results page for the query "Seattle". The browser's address bar shows the URL `http://search.live.com/results.aspx?q=Seattle&mkt=en-gb&FORM=LVCP`. The page features a search bar with "Seattle" entered and a "Search" button. Below the search bar, there are several search results, including "Seattle Flights", "Visiting Seattle?", and "Seattle.gov".

Overlaid on the screenshot are several mathematical annotations:

- At the top right, the formula  $b_1 \cdot p_1 \geq b_2 \cdot p_2$  is displayed, with arrows pointing to the search results for "Seattle Flights" and "Seattle Times".
- In the middle right, the formula  $c_i = b_{i+1} \cdot \frac{p_{i+1}}{p_i}$  is shown, with arrows pointing to the "Seattle Flights" result and the "Seattle Times" result.

The search results are presented in a table-like format with colored boxes:

Seattle Flights - www...	\$1.00	* 10%	=\$0.10	\$0.80
Visiting Seattle? - S...	\$2.00	* 4%	=\$0.08	\$1.25
seattle - www.gawwk.o...	\$0.10	* 50%	=\$0.05	\$0.05

Related searches listed on the right include: Seattle Weather, Seattle Times, Seattle Hotels, Craigslist Seattle, Seattle Washington, Seattle Mariners, Craigs List Seattle, and Seattle Seahawks.

Sponsored sites include: Seattle Washington Rates, Visiting Seattle & Need a Hotel?, Seattle Washington Hotel Bargains!, and www.nextag.com/hotels.

# AdPredictor Technology Transfer





# SQL Schema Generator

- **Code size:** 500 LOC
- **Project size:** 1 file
- **Development time:** 2 weeks
  
- **Features**
  - Code defines the schema (unlike LINQ)!
  - High-performance insertion via computed bulk-insertion with automated key propagation
  - Code sample is now part of the F# distribution

# Strong Typing and SQL Datastores

```
/// A single page-view  
type PageView =
```

```
{  
    ClientDateTime : DateTime
```

```
/// Create the SQL schema
```

```
let schema = bulkBuild ("cpidssdm18", "Cambridge", "June10")
```

```
/// Try to open the CSV file and read it pageview by pageview
```

```
File.OpenTextReader "HourlyRelevanceFeed.csv"
```

```
|> Seq.map (fun s -> s.Split [|', '|])
```

```
|> Seq.chunkBy (fun xs -> xs.[0])
```

```
|> Seq.iteri (fun i (rguid,xss) ->
```

```
    /// Write the current in-memory bulk to the Sql database
```

```
    if i % 10000 = 0 then
```

```
        schema.Flush ()
```

```
    /// Get the strongly typed object from the list of CSV file lines
```

```
    let pageView = PageView.Parse xss
```

```
    /// Insert it
```

```
    pageView |> schema.Insert
```

```
/// One final flush
```

```
schema.Flush ()
```

```
LocationMetroArea : int
```

```
/// Different types of media
```

```
type MediumType =
```

```
| PaidSearch
```

```
| ContextualSearch
```

```
OrderItemAdPlacement
```

```
OrderItemAdPlacement
```

```
OrderItemId : int
```

```
OrderItemId : int
```

```
CampDayId : int16
```

```
CampHourNum : byte
```

```
ProductId : ProductType
```

```
MatchType : MatchType
```

```
AdLayoutId : AdLayout
```

```
RelativePosition : byte
```

```
DeliveryEngineBank : int16
```

```
ActualBid : int
```

```
Check : int16
```

```
MatchScore : int
```

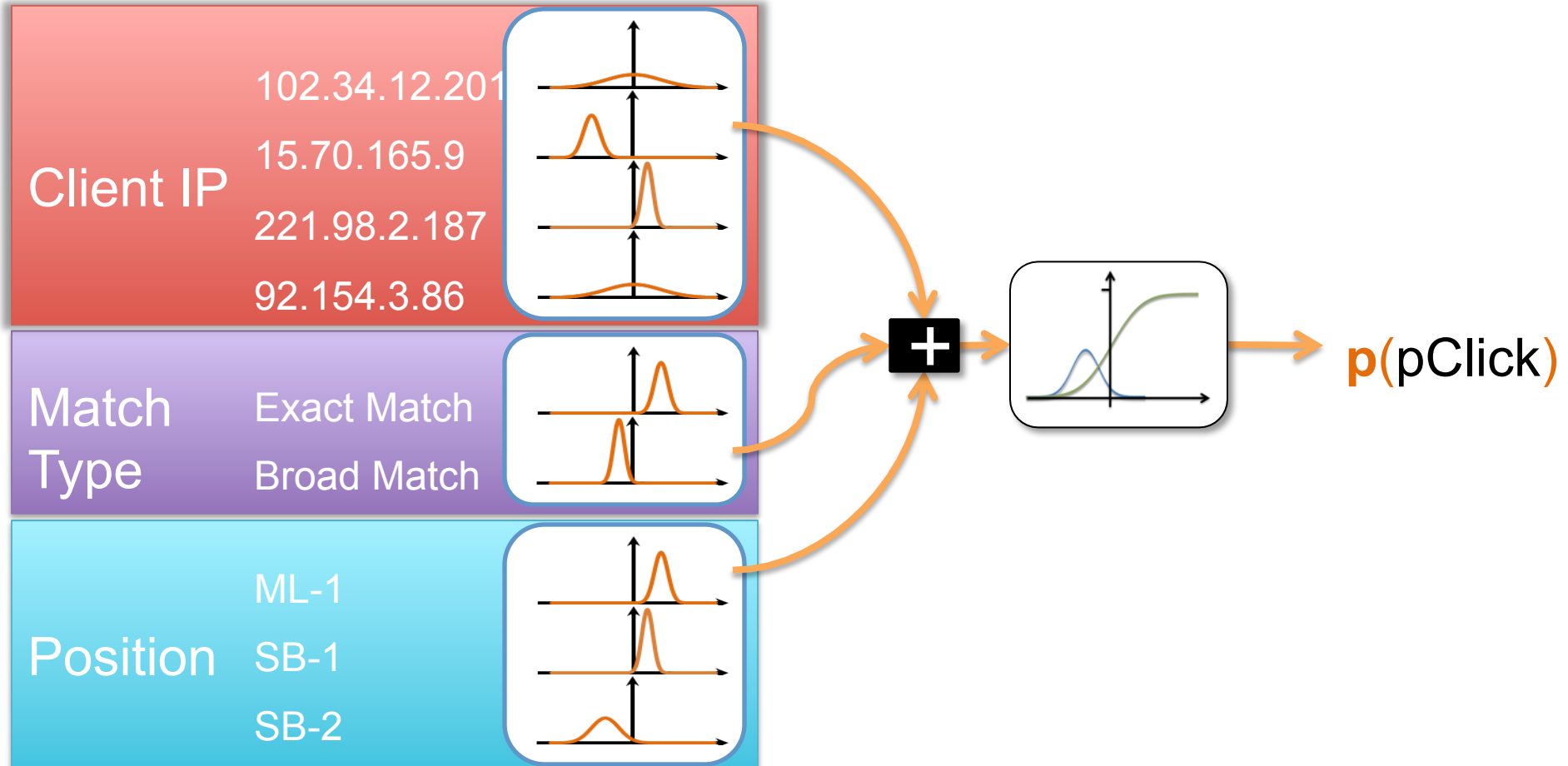
```
ImpressionCnt : int
```

```
ClickCnt : int
```

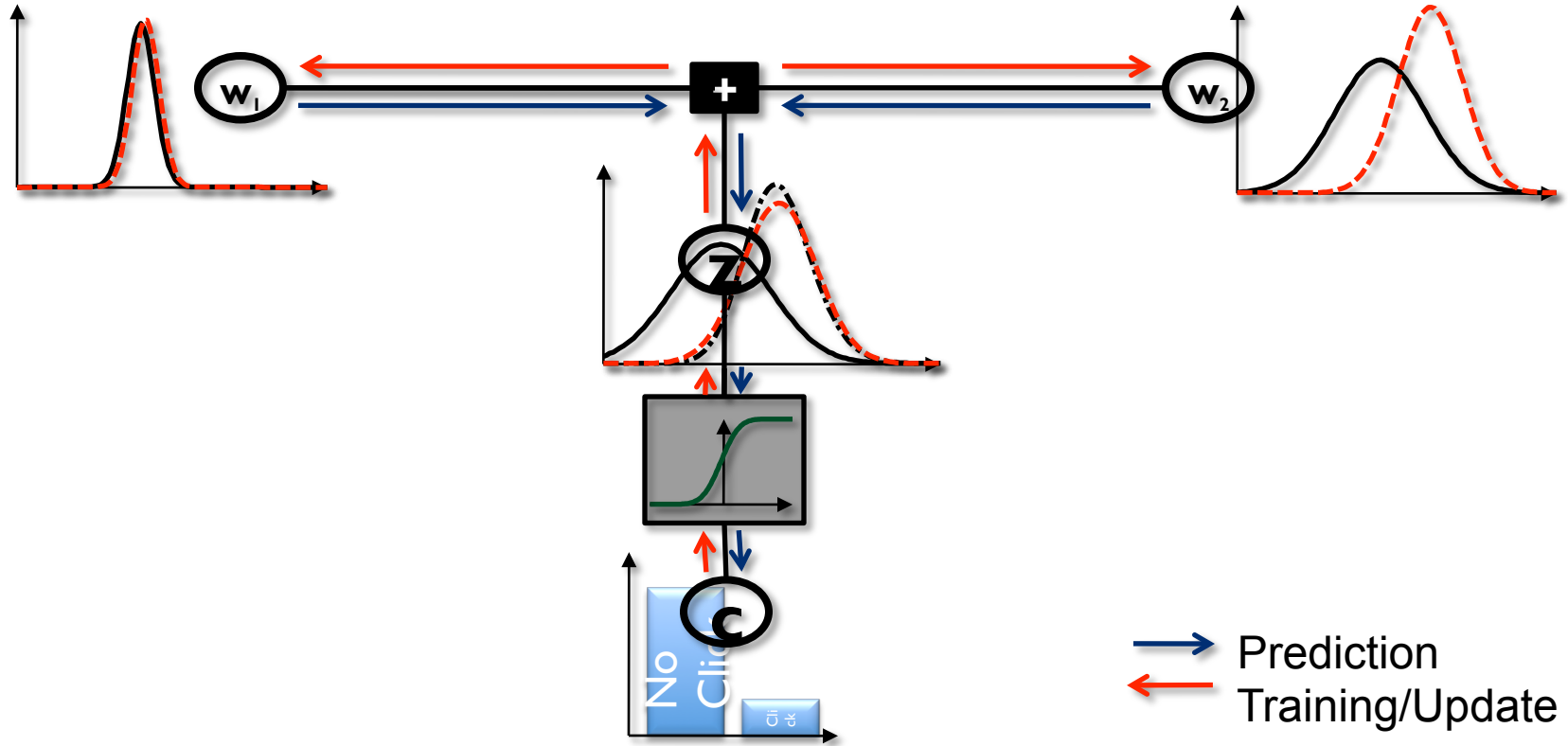
```
ConversionCnt : int
```

```
TotalCost : int
```

# Uncertainty: Bayesian Probabilities



# Training Algorithm in Action

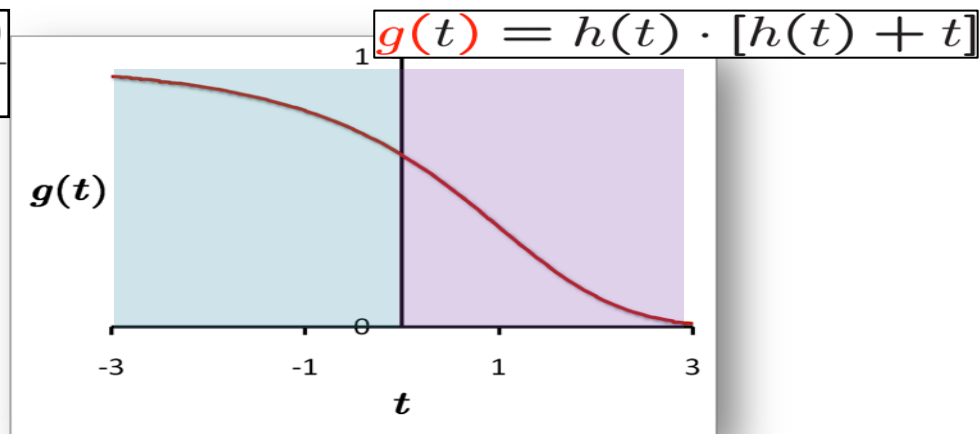
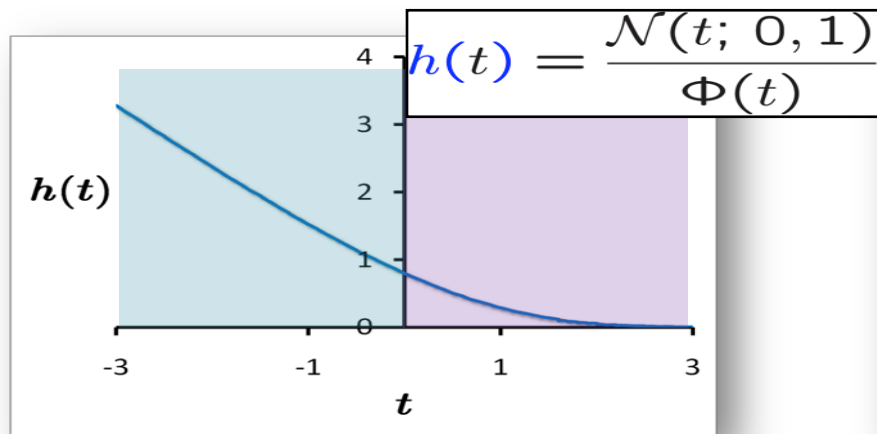




# Inference: An Optimization View

$$\mu_i \leftarrow \mu_i + \frac{\sigma_i^2}{s} \cdot h \left[ \frac{\sum_{j=1}^d \mu_j}{s} \right] \quad \sigma_i^2 \leftarrow \sigma_i^2 \left( 1 - \frac{\sigma_i^2}{s^2} \cdot g \left[ \frac{\sum_{j=1}^d \mu_j}{s} \right] \right)$$

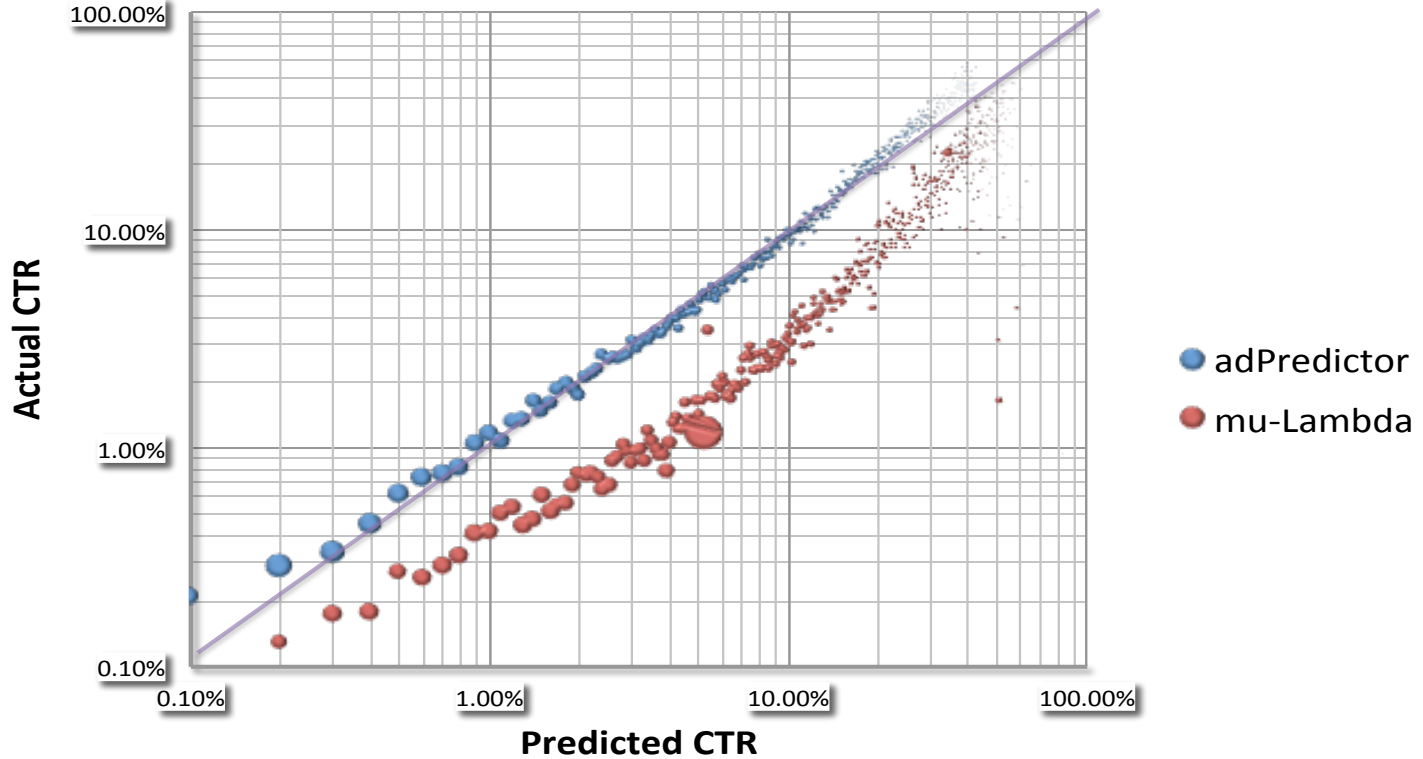
$$s^2 = \beta^2 + \sum_{j=1}^d \sigma_j^2$$



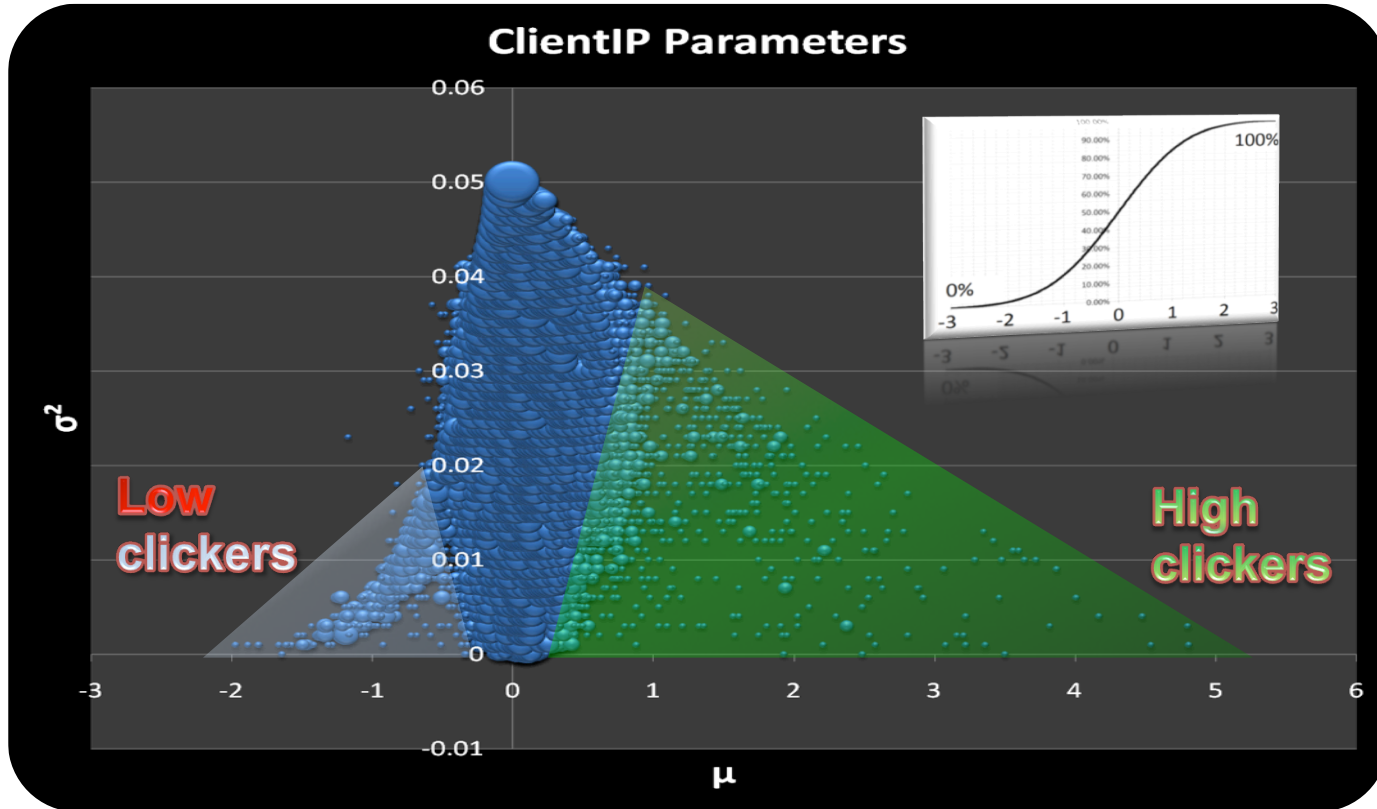
# AdPredictor Technology Transfer



# Offline Evaluation



# Client IP: Mean & Variance



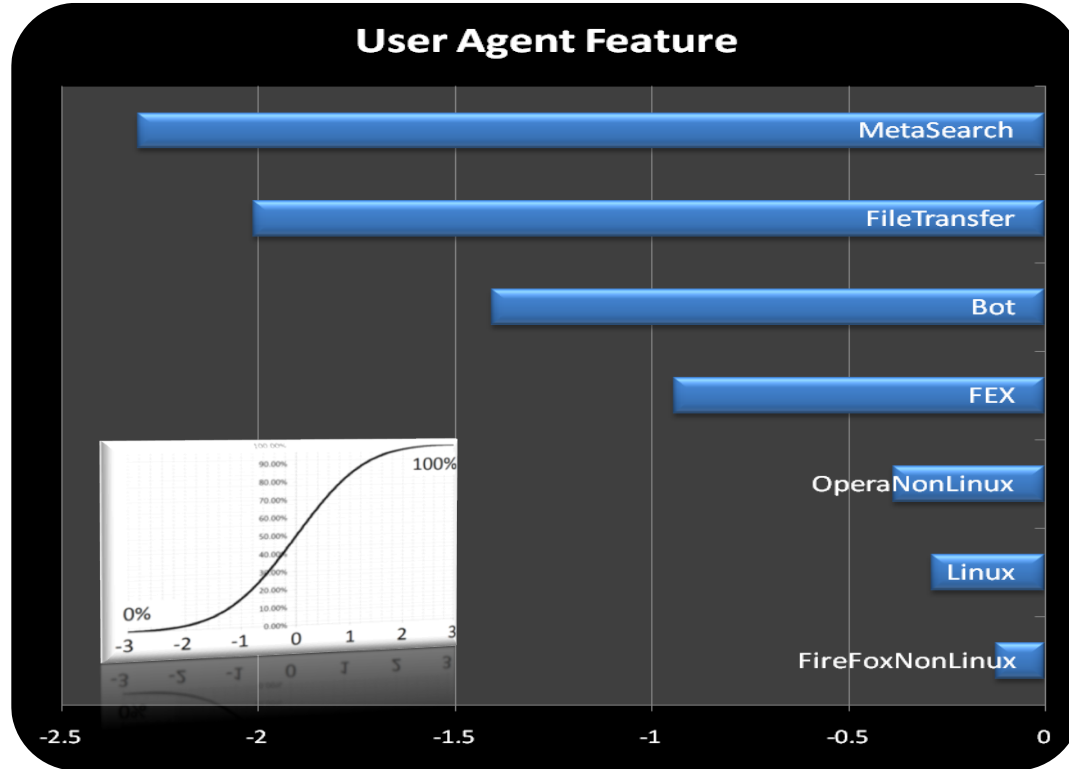
**Low  
clickers**

**High  
clickers**

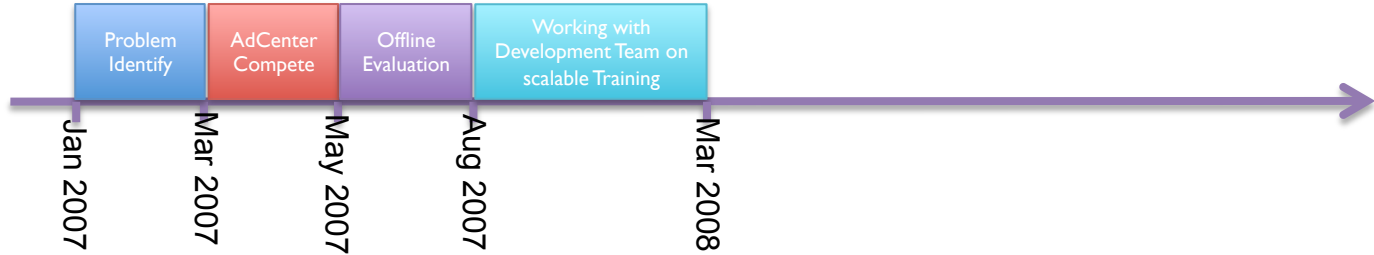
$\mu$

$\sigma^2$

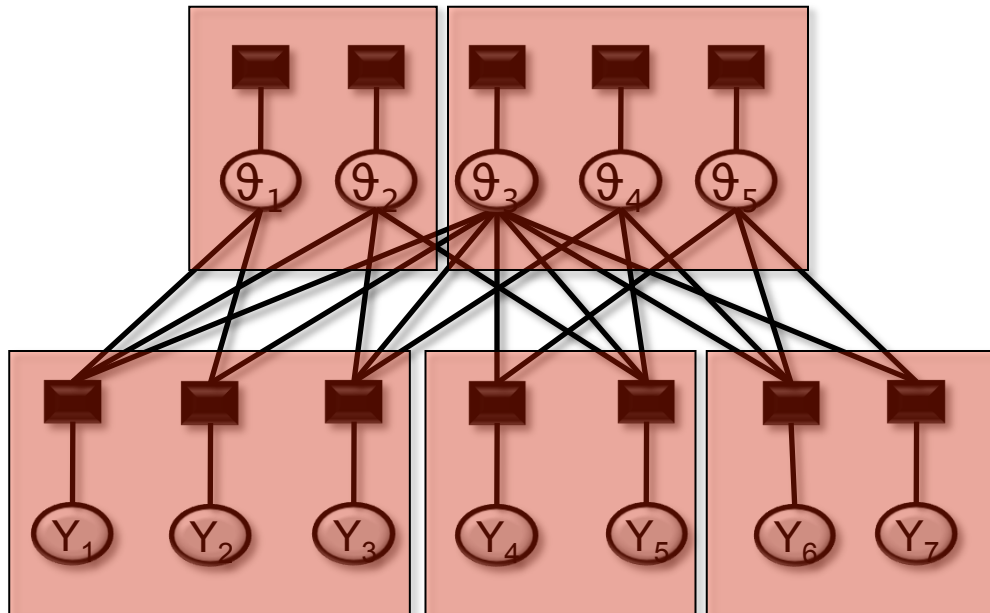
# UserAgent: Mean Posterior Effects



# AdPredictor Technology Transfer



# Distributed Conditional Models



**Belief Store**  
("Memory")

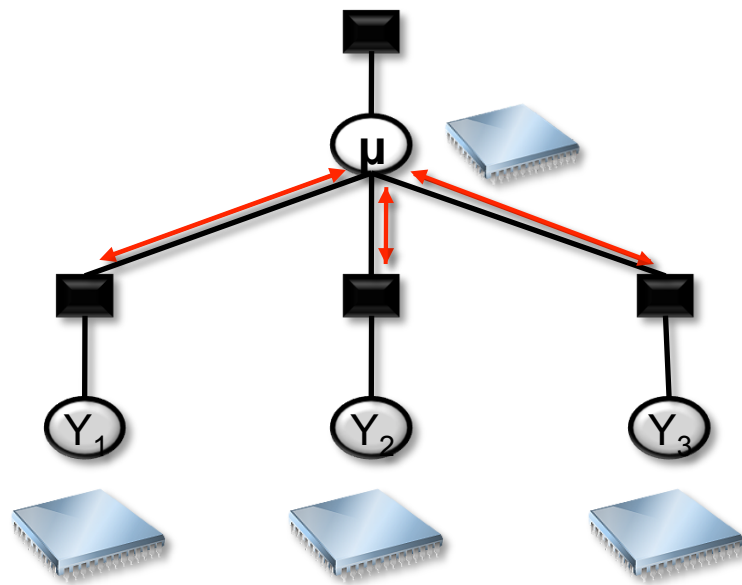
**Message Passing**  
("Communicate")

**Data Messages**  
("Compute")

# Relation to Map-Reduce

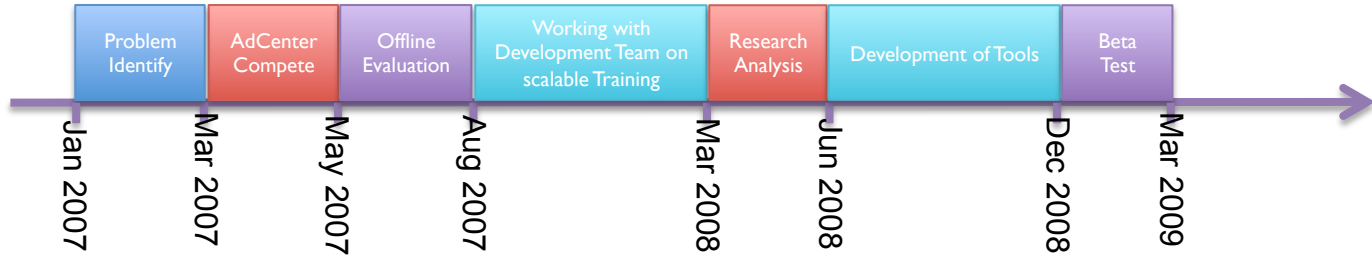
- **Map-Reduce**
  - Map: Data nodes compute messages  $m_{F_k \rightarrow \theta}$  from data  $y_i$  and  $m_{\theta \rightarrow F_k}$
  - Reduce: Combine messages  $m_{F_k \rightarrow \mu}$  into  $p(\theta)$  by multiplication
  - Vanilla MR is a single pass only!
- **Caveats:**
  - Approximate data factors need all incoming message  $m_{F_k \rightarrow \theta}$ !
  - Each machine needs to be able to store the belief over  $\theta$

$$p(\theta | \mathbf{x}, \mathbf{y}) \propto \prod_k f_k(Y_k | \theta, X_k) \cdot p(\theta)$$



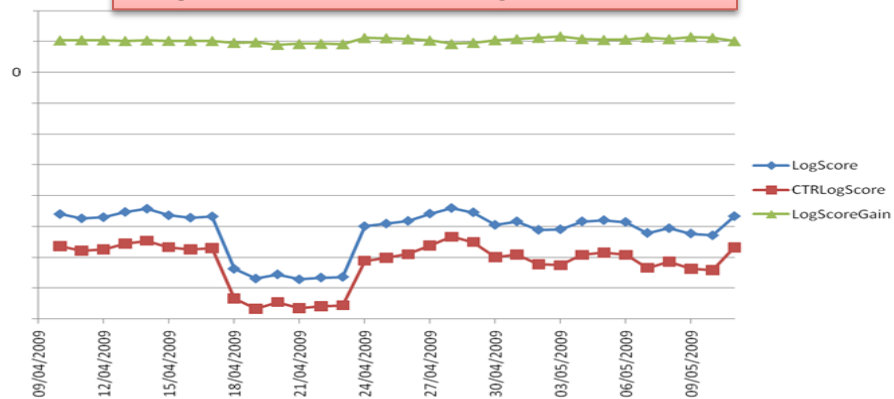


# AdPredictor Technology Transfer

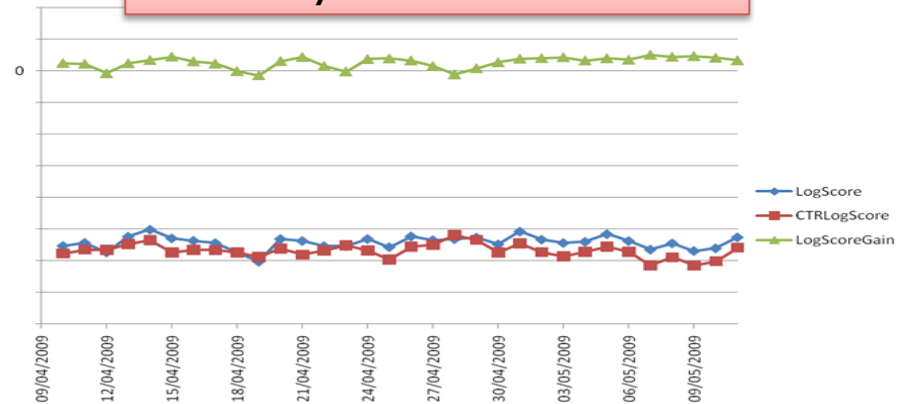


# Online Metrics

## Bayesian Probit Regression



## Naive Bayes



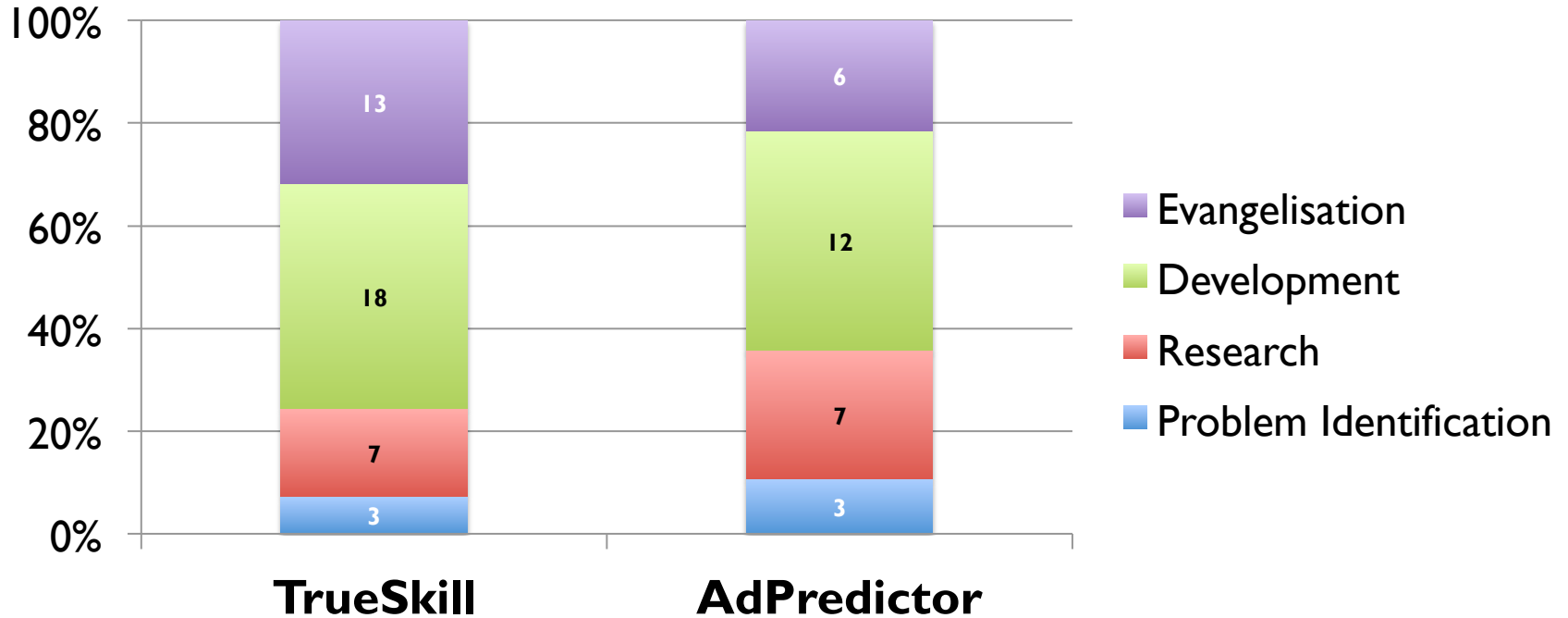
# AdPredictor Technology Transfer



- **Lessons Learned:**

1. Pure research takes a short amount of time
2. Development takes much longer than planned
3. Metrics are important and part of the transfer
4. Develop for scale from Day 1

# Technology Transfer in Numbers



# Overview

- Background
- Technology Transfer Case Studies
  - TrueSkill: Gamer Rating and Matchmaking
  - Click-Through Rate Prediction in Online Advertising
- **Technology Transfer Lessons**
  - Process
  - Technical

# Process Lessons

- Identify the problem
- Identify the customer
- Be the customer!

**Problem vs. Technique**

- Do not study the existing literature first!
- Don't be afraid to be wrong

**Try Out and Then Study**

- Reduce risk!
- Understand the decision and engineering process
- Respect timelines

**Under-Promise and Over-Deliver**

- Write code!
- Work with developers – not executives!

**Coding**

- Simplify to its bare minimum.
- Develop tools and help adoption with customers

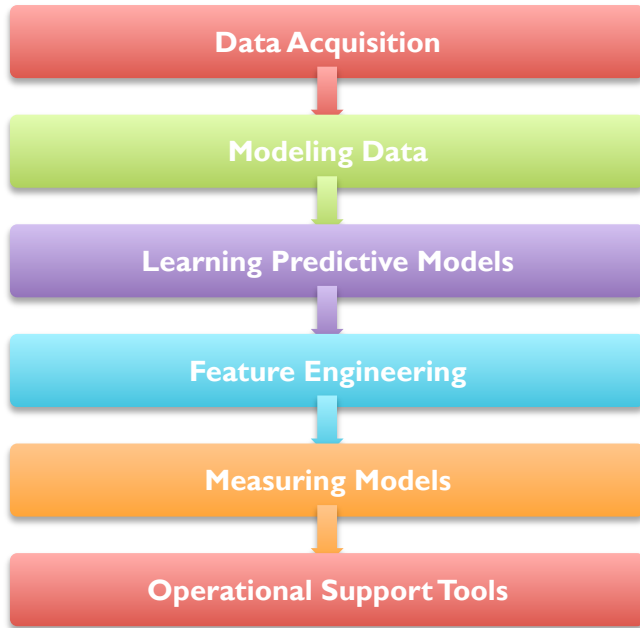
**Simplicity**

- Be in for the long haul – years is normal
- Do not aim for a quick win!

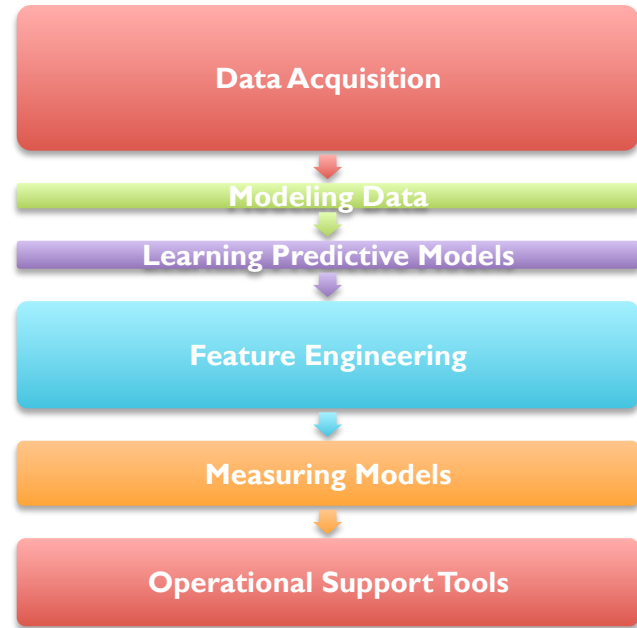
**Long Run**

# Process Lessons: Pictures

## Steps



## Time Allocation



# Technical Lessons: Practical Problems

## Ideal

- Business case well defined
- Data pipeline established
- Training set given
- Business metric/loss given
- Meaning of data fields fixed
- Breakthrough impact through ML algorithm

## Reality

- Business case unclear
- Irregular data-file drops
- No training set
- Unclear measure impact
- Missing & inconsistent data
- ML algorithm leads to single digit improvement at best



# How to Pick a Practical ML Problem

- **Key Questions:**
  1. **Data:** Will we have sufficient and ongoing data?
  2. **Complexity:** Can a simple rule do as well?
  3. **Customer Experience:** How will the customer see the prediction/summarization? What will it impact?
  4. **Economics:** What's the cost and benefit of a single prediction?

# Example of Practical ML Problems

## Good

- Click-Through-Rate Prediction
- Demand Forecasting
- Named Entity Extraction
- Fraud Prediction

## Bad

- **Data:** Prediction of mushroom/flower types
- **Complexity:** Learning to predict imputed data
- **Customer Experience:** Predictive menus
- **Economics:** Complex, non-linear models for advertising

# Conclusion

- Technology Transfer is Highly Rewarding!
- Practical problems start new research directions!
- Graphical models are a very powerful language:
  - Modeling (Bayes Nets)
  - Algorithm development (Sum-Product)
  - Highly modular (Local Factors)
  - (Relatively) easy to teach (Pictorial)
- Machine Learning = “Statistic of Big Data”?