IBM Research: Mobile, Solutions, and Mathematical Sciences

# Data Analytics in Industry Research:
# A Personal Perspective*

Naoki Abe
Senior Manager, Data Analytics
Mathematical Sciences and Analytics
IBM T. J. Watson Research Center

*Based on joint work with H. Mamitsuka, A. Nakamura, H. Li, J. Takeuchi, E. Pednault, P. Melville, S. Rosset, Y. Liu, A. Lozano, R. Luss P. Olsen, E. Yang, K. Ramamurthy, M. Kshirsagar, et al

# Agenda

- NEC days (1990's)
  - Stochastic tree grammars and Genetic Information Processing
  - On-line binary relation learning and Natural Language Processing
  - ...

- IBM days (2000's - present)
  - Smarter Government
    - Markov Decision Process and Tax Collections
  - Smarter Planet , Enterprise & Cloud
    - Temporal Causal Modeling and Applications
  - Smarter Agriculture
    - Heterogeneous Data Analytics and Accelerated Plant Breeding

- Summary and Retrospective

IBM Research: Mobile, Solutions, and Mathematical Sciences

# NEC and History of IBIS...

- Founding members of "AI Basic Research Group" (Katsuhiro Nakamura, Kenji Yamanishi, Junichi Takeuchi) had a great idea

- "IBIS" = Endangered Species





IBIS'98
1998年情報論的学習理論ワークショップ開催案内
IBIS'98
(1998 Workshop on Information-Based Induction Sciences)

● 主催
　情報理論とその応用学会(SITA)
　協賛
　電子情報通信学会情報理論研究専門委員会
　人工知能学会
● 開催期間:
　1998年7月11日(土)13:20 --7月12日(日)17:00
● 開催場所:
　専修大学箱根セミナーハウス
　神奈川県足柄下郡箱根町元箱根字大芝103

● 問い合わせ先

　●実行世話役
　●ワークショップ全般
　山西　健司 (tel: [電話番号削除] fax: [FAX番号削除], email: [メールアドレス番号削除])
　中村　勝洋 (tel: [電話番号削除] fax: [FAX番号削除], email: [メールアドレス番号削除])
　〒216-8555 川崎市宮前区宮崎4-1-1, NEC C&Cメディア研究所

　● 予稿集関連
　竹内　純一 (tel: [電話番号削除] fax: [FAX番号削除], email: [メールアドレス番号削除])
　〒216-8555 川崎市宮前区宮崎4-1-1, NEC C&Cメディア研究所

　● 専修大学箱根セミナーハウス関連(一般質問)
　佐藤　創 (tel: [電話番号削除] email: [メールアドレス番号削除])
　川崎市多摩区東三田2-1-1,専修大学経営学部情報管理学科
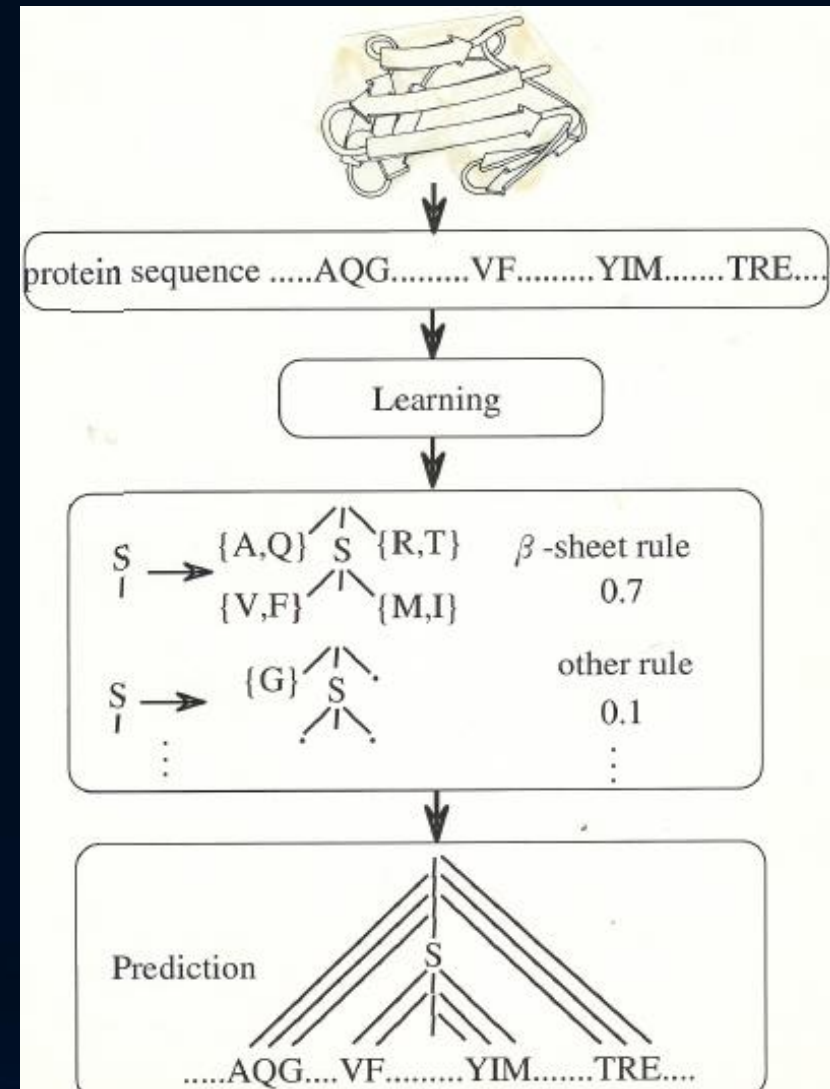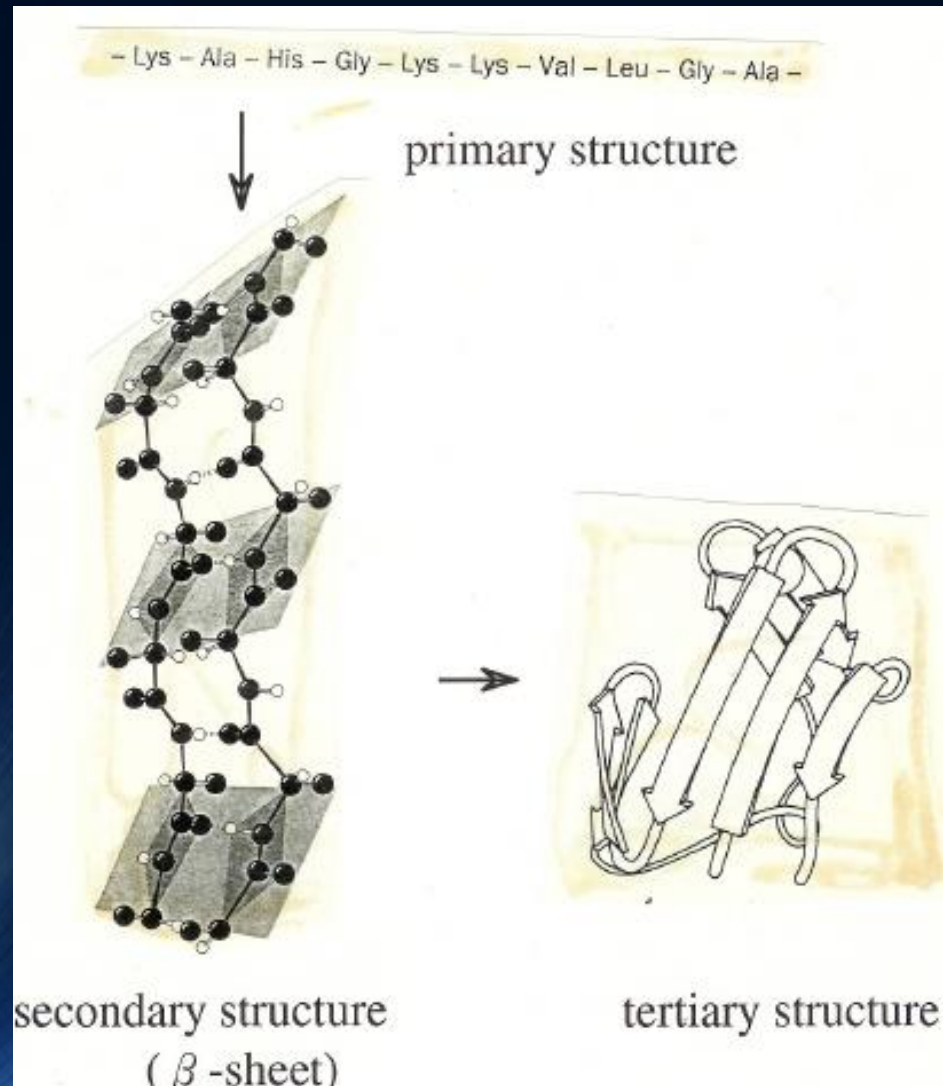
　● 情報理論とその応用学会関連(一般質問)
　企画理事　山口　和彦(電通大)/岩村　恵一(キヤノン)

# NEC Days: 1990's

- ***Stochastic Tree Grammars and Genetic Information Processing*** (H. Mamitsuka)

- On-line Binary Relation Learning and Natural Language Processing (with A. Nakamura, H. Li)

- On-line Active Learning and Rational Choice Theory (with J. Takeuchi, S. Amari)

- On-line Active Learning and Internet Advertisement Optimization (with A. Nakamura, et al)

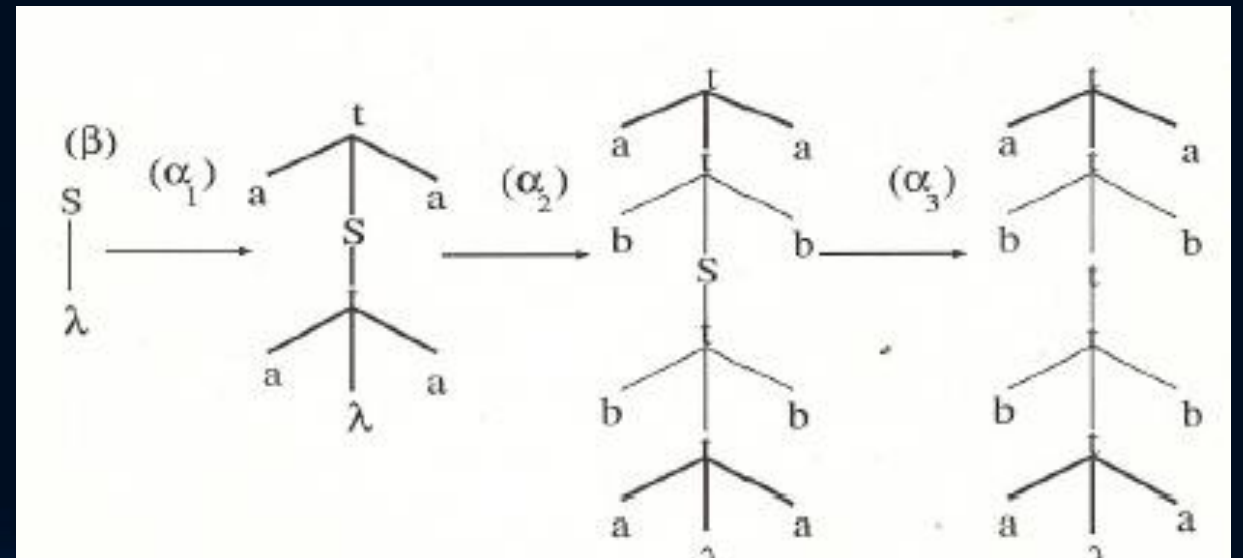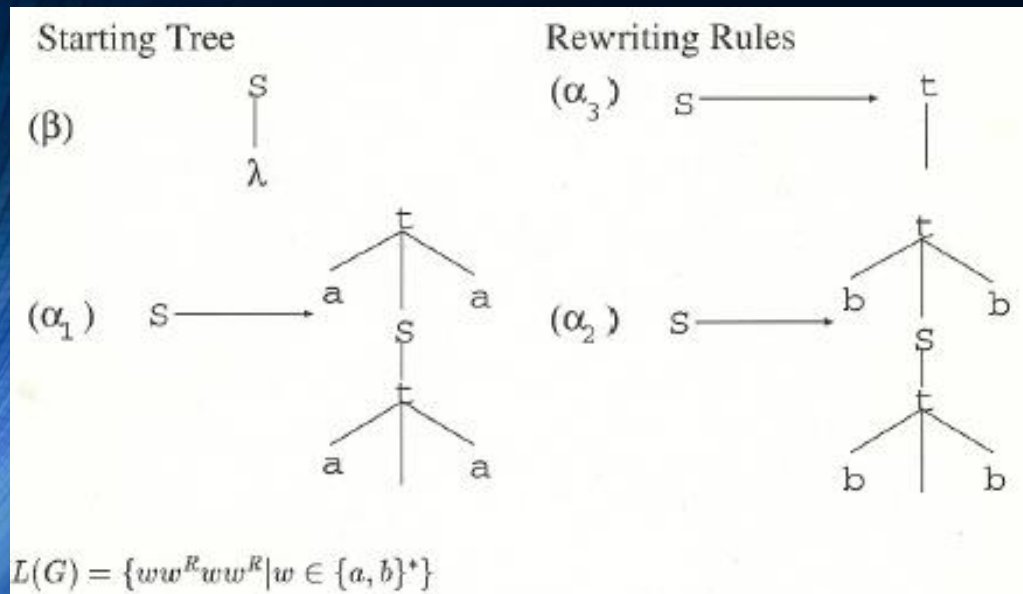- Active Learning and Immunological Experimental Design (H. Mamitsuka, K. Udaka, et al)
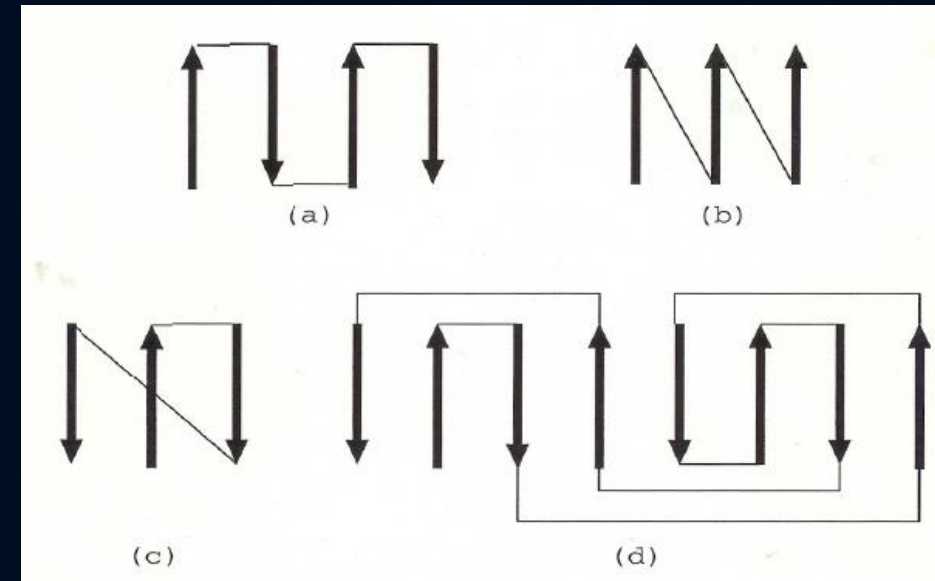
- …

# Predicting Protein Secondary Structure (Beta-sheet Regions) with Stochastic Tree Grammars
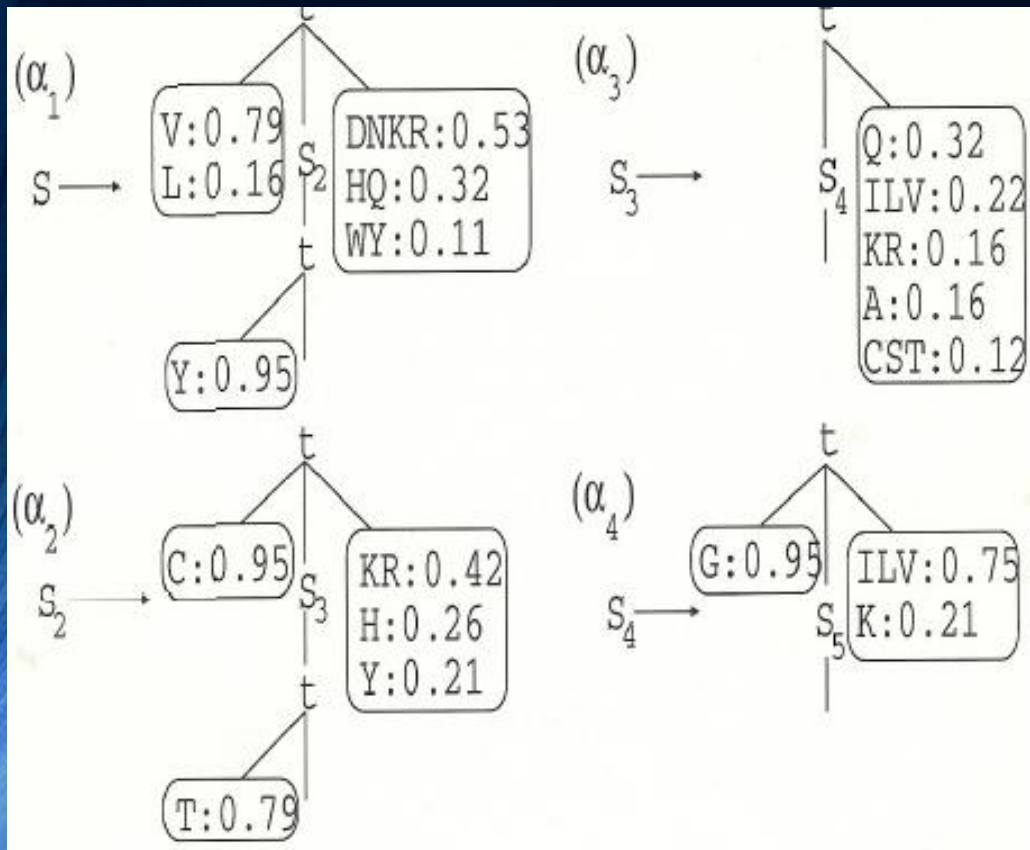
H. Mamitsuka & N. Abe

# Beta-sheet Patterns and Tree Grammars



- Beta sheet patterns exhibit so-called "unbounded dependencies"

- Ranked Node Rewriting Grammar RNRG inspired by (and generalizes) Tree Adjoining Grammars (Joshi et al)
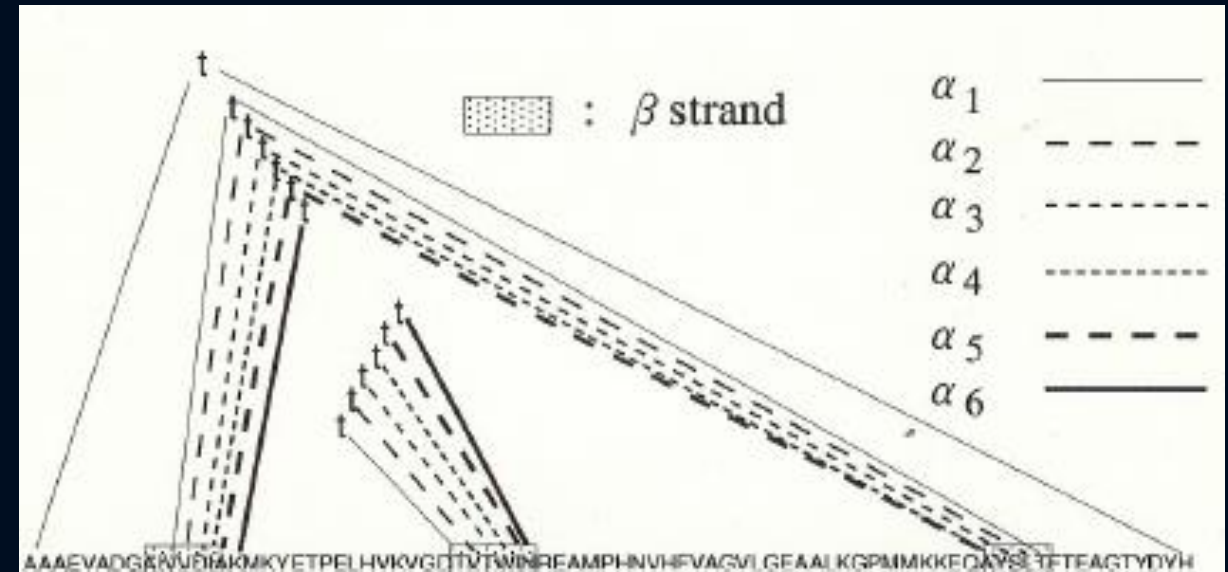


$$L(G) = \{ww^Rww^R | w \in \{a, b\}^*\}$$

# Sample results of predicting beta-strand in a Toxyn

- Learned grammar

- Most likely parse

# NEC Days: 1990's

- Stochastic Tree Grammars and Genetic Information Processing (H. Mamitsuka)

- ***On-line Binary Relation Learning and Natural Language Processing*** (with A. Nakamura, H. Li)

- On-line Active Learning and Rational Choice Theory (with J. Takeuchi)

- On-line Active Learning and Internet Advertisement Optimization (with A. Nakamura)

- Active Learning and Immunological Experimental Design (H. Mamitsuka)

- …

# On-line binary relation learning and word clustering

A. Nakamura, H. Li & N. Abe

- Sub-categorization in NLP: naturally viewed as clustering via density estimation

- An Alternative Formulation: On-line Learning of Binary Relations





On-line Learning: Minimize number of prediction errors!

# Two-dimensional Weighted Majority Algorithm Out-Performs Heuristic Methods



WMP2 (Conservative Version)



Upper Bound on Worst Case Number of Mistakes

## Theorem

Algorithm WMP2 makes at most

$$\frac{1}{k+l}\left(kl(m+n)+(ln+km)\sqrt{2(m+n)\log\frac{kl(m+n)}{ln+km}}\right)$$

mistakes when learning a (k,l)-binary relation



WMP2

2-dimensional WMPs

1-dimensional WMPs

- 2-dimensional WMP algorithms outperform 1-dimensional ones

- Theoretically sound algorithm (WMP2) outperforms heuristic ones…

# IBM Days (2000's – present)

- ***Smarter Government***
  - ***Markov Decision Process and Tax Collections***

- Smarter Planet , Enterprise & Cloud
  - Temporal Causal Modeling and Applications

- Smarter Agriculture
  - Heterogeneous Data Analytics and Accelerated Plant Breeding

IBM Research: Mobile, Solutions, and Mathematical Sciences

# Smarter Government

- Tax Collections Optimization System (TACOS) video
  - https://www.youtube.com/watch?v=VGKp13APsBg

# Tax Collections Optimization for NYS: Technical Challenge

- Tax collections process is a complex process involving various legal/business constraints

- Most existing approaches rely on rigid, manual rules, including NYS legacy system

- Goal: take this rigid procedure apart, leaving fragments of it intact wherever necessary, and automatically configure the rest, based on analytics and optimization



A Fragment of Legacy Process

# Tax Collections Optimization for NYS: Technical Challenge

- Tax collections process is a complex process involving various legal/business constraints
- Most existing approaches rely on rigid, manual rules, including NYS legacy system
- Goal: take this rigid procedure apart, leaving fragments of it intact wherever necessary, and automatically configure the rest, based on analytics and optimization
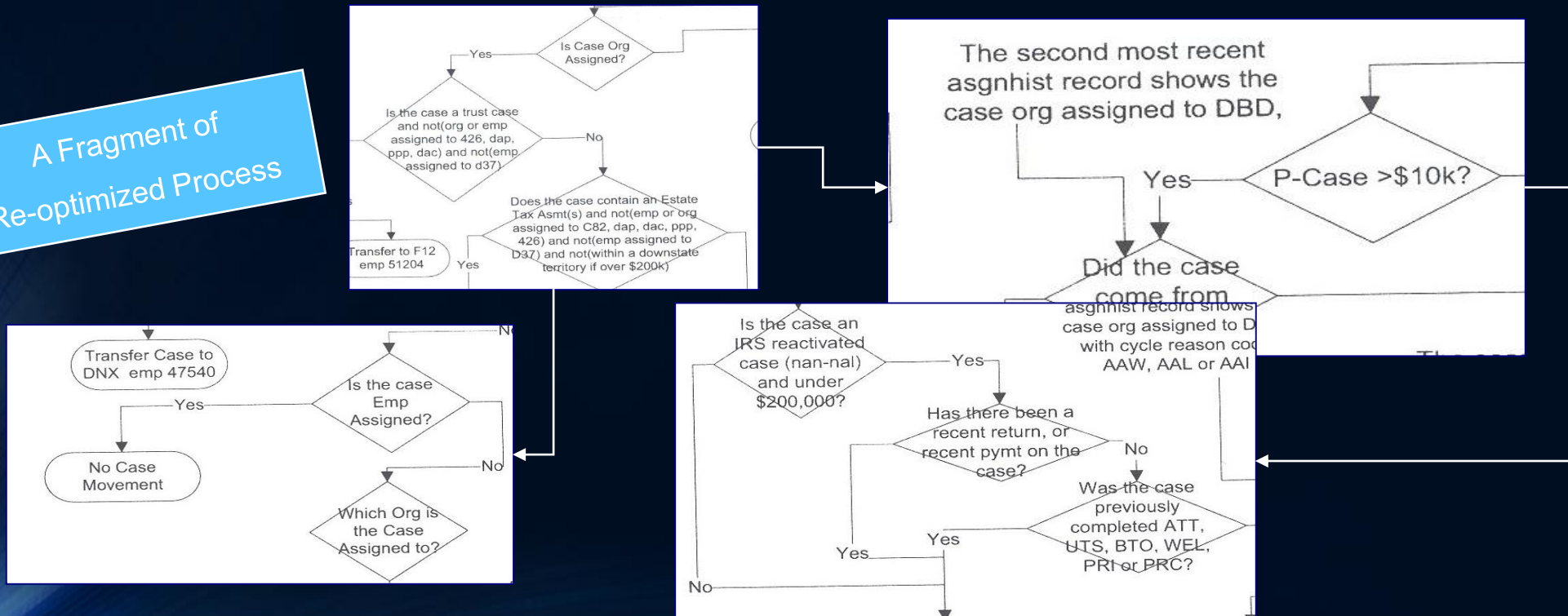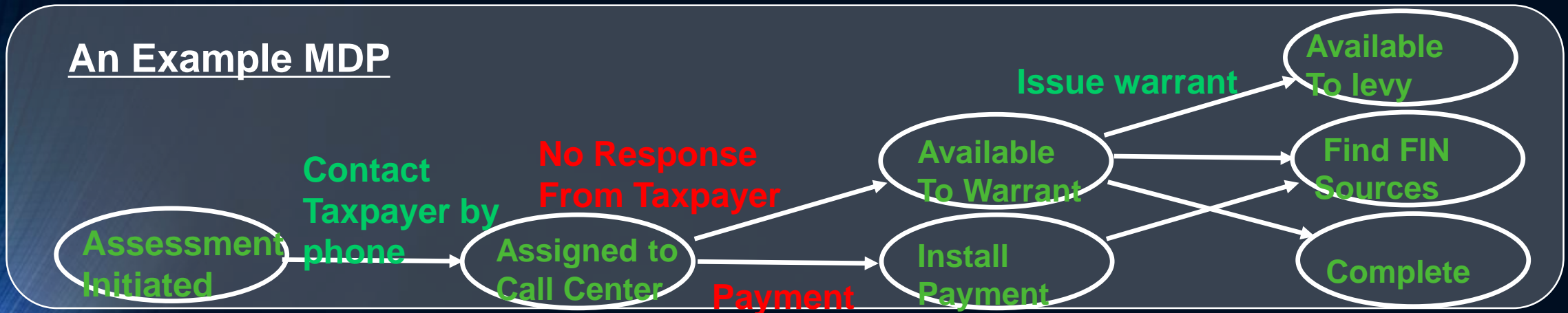


A Fragment of Re-optimized Process

# Tax Collections Optimization for NYS: The Framework

- **Markov Decision Process (MDP) formulation provides an advanced framework for modeling tax collection process**
    - "States", $s$, summarize information on a taxpayer's stage in collection process
    - "Action", $a$, is a collection action (e.g. phone call, warrant, levy)
    - "Reward", $r$, is the tax collected for the taxpayer in question

**An Example MDP**

Issue warrant

Available To levy

Find FIN Sources

No Response From Taxpayer

Available To Warrant

Contact Taxpayer by phone

Assessment Initiated

Assigned to Call Center

Install Payment

Complete

Payment

- The goal in MDP is formulated as outputting a policy which maps TP's states to collection actions so as to maximize the long term cumulative rewards
- **Constrained MDP** requires additionally that output policy belongs to a constrained class adhering to certain constraints

# Methodology: Constrained Reinforcement Learning (C-RL)

- **Business requirement to customize collections actions depending on detailed taxpayer characteristics**
  - Use of a high-dimensional state space necessary
- **With high-dimensional state space**
  - Estimating the structure of MDP is extremely challenging
  - Reinforcement Learning (RL) solves MDP with access only to data (not MDP itself)
  - We develop constrained-RL (C-RL) methods for high dimensional state space

# Tax Collections Optimization for NYS: Example Segment

**Segment Definition**

    **state = Call-Center-Not-Warranted**

    **and  tax_pd_lst_yr < $X**

    **and 1 <= num_st_ff_cllct_asmts**

    **and 1 <= num_non_rstrctd_fin_srcs**

    **and 1 <= st_inactv_ind**
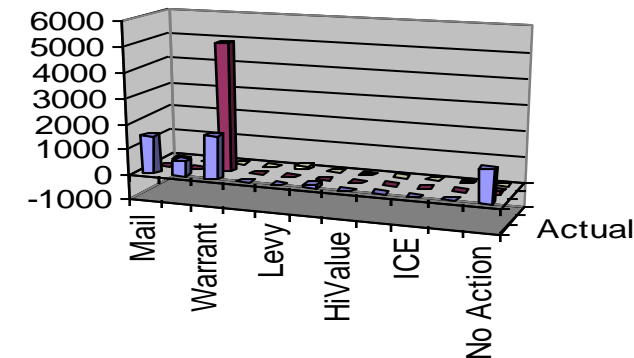
    **and  num_pymnts_snc_lst_actn < 1**

- **Interpretation**
  - Case is in Call Center and has not been warranted
  - There is at least one non restricted fin source identified
  - Sales tax inactive indicator is on
  - There was no payment in the last period
  - Tax paid last year is less than 1000 dollars
- **Create warrant recommended**

Segment size 1.0%

**Action Distributions Segment 212**



| | Mail | Phone | Warrant | FS IE | Levy | DO | HiValue | CVS | ICE | Field Visit | No Action |
|---|---|---|---|---|---|---|---|---|---|---|---|
| ▦ Actual | 1522 | 612 | 1702 | 0 | 3 | 126 | 2 | 0 | 0 | 3 | 1293 |
| ▦ Allocated | 0 | 0 | 5103 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 152 |
| ☐ Coefficients | -0.208 | 0.821 | 1.934 | 0 | 130.22 | 0.662 | -0.109 | 0 | 0 | -0.139 | 0 |

# Deployment Results: Tax Collection Optimization for NYS

- 2010 realized an 8.22% revenue increase over 2009 - $+83M

| Statistics | |
|---|---|
| Average age of cases | 9.3% decrease |
| Dollar per staff day for field agents | 15% Increase |
| Collection by field staff | 12% increase |
| Dollar per warrant | 22% increase |
| Dollar per levy | 11% increase |
| Number of warrants filed | 9% decrease |
| Number of levies served | 3% decrease |
| 35,000 less taxpayers had serious actions taken against them | |

*Many taxpayers that would have had action taken in the past had no action taken yet debt was still collected*

Edelman Presentation NYS DTF Final updated 4-11-2011

Recognition in Academia:

- **Awarded "best government & industry paper" at KDD 2010**

- **Finalist for Edelman prize of INFORMS 2011**

# IBM Days (2000's – present)

- Smarter Government

  - Markov Decision Process and Tax Collections

- ***Smarter Planet , Enterprise & Cloud***

  - ***Temporal Causal Modeling and Applications***

- Smarter Agriculture

  - Heterogeneous Data Analytics and Accelerated Plant Breeding

IBM Research: Mobile, Solutions, and Mathematical Sciences

# Smarter Planet (and Enterprise and Cloud…)

• "The most ridiculous thing I've ever heard…"



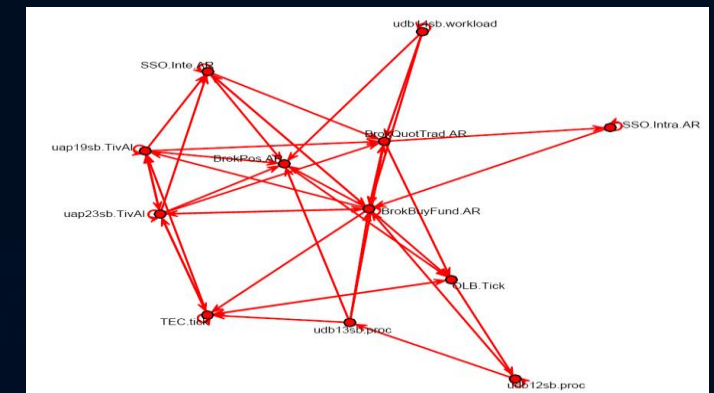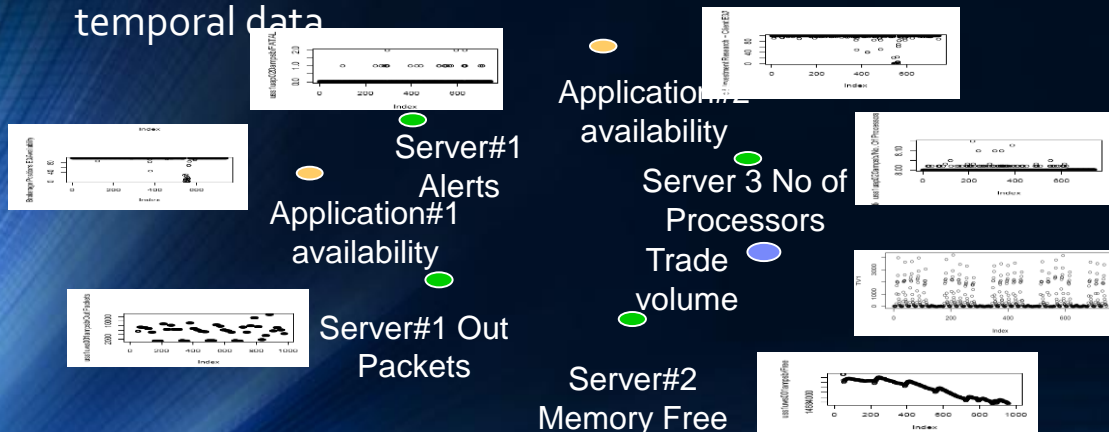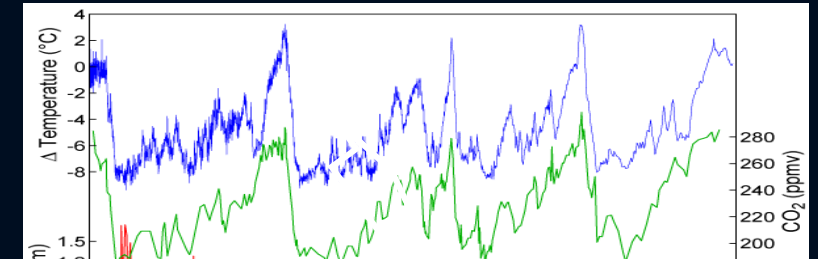Source: "An Inconvenient Truth" the movie

# Temporal Causal Modeling by Graphical Granger Modeling

- Granger causality

  - First introduced by the Nobel prize winning economist, Clive Granger

- Definition: a time series x is said to "Granger cause" another time series y, if and only if regressing for y in terms of both past values of y and x is statically significantly better than that of regressing in terms of past values of y only

$$y_t \approx A \cdot y_{\overrightarrow{t-1}} + B \cdot x_{\overrightarrow{t-1}} \qquad (1)$$

$$y_t \approx A \cdot y_{\overrightarrow{t-1}} \qquad (2)$$

- Combination of Granger Causality and cutting-edge graphical modeling techniques provides efficient and effective methodology for graphical causal modeling of temporal data



Application#2 availability

Server#1 Alerts

Server 3 No of Processors Trade volume

Application#1 availability

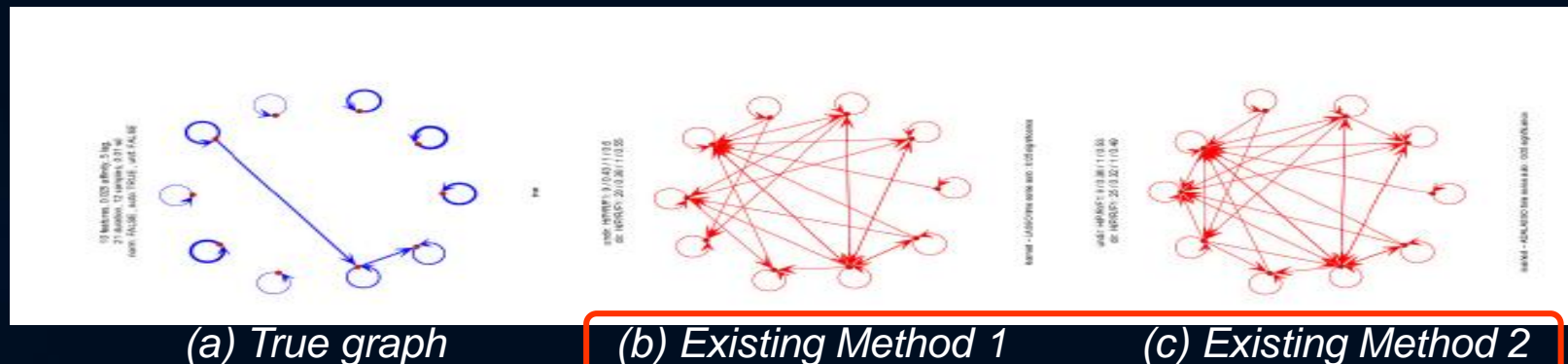Server#1 Out Packets

Server#2 Memory Free
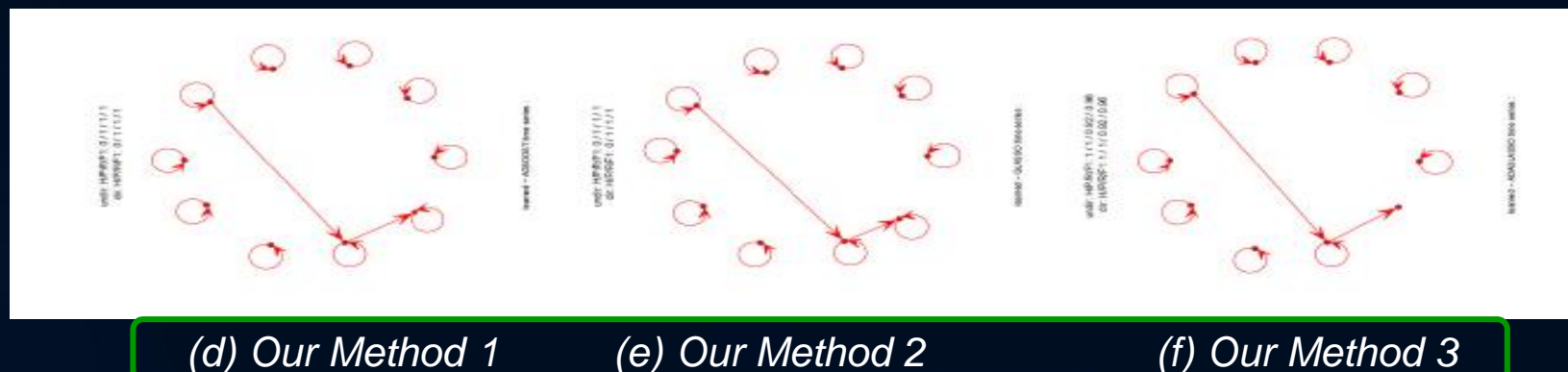
# Novel Graphical Granger Modeling Methods

- Our methodology leverages temporal constraints in graphical Granger modeling by treating lagged variables of the same feature as a group, and invokes ***"structured sparse modeling"*** technique

*Example Outputs*

**Existing methods**



(a) True graph      (b) Existing Method 1      (c) Existing Method 2

**Our methods**



(d) Our Method 1      (e) Our Method 2      (f) Our Method 3

*Accuracy Comparison*

| Method | Existing 1 | Existing 2 | Our 1 | Our 2 | Our 3 |
|---|---|---|---|---|---|
| Accuracy(%) | 62 | 65 | 92 | 87 | 91 |

# Climate change attribution by spatial-temporal causal modeling

- Identify causal relationships between natural/anthropogenic forcing factors and climate variables, based on spatio-temporal observations.



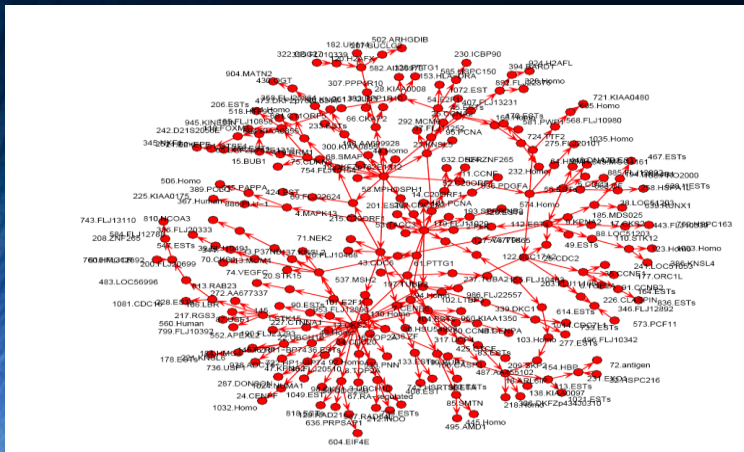| Variables (Variable group) | Type | Source |
|---|---|---|
| Methane ($CH_4$)<br>Carbon-Dioxide ($CO_2$)<br>Hydrogen ($H_2$)<br>Carbon-Monoxide (CO) | Greenhouse Gases | NOAA |
| UV (AER) | Aerosol Index | NASA |
| Temperature (TMP)<br>Temp Range (TMP)<br>Temp Min (TMP)<br>Temp Max (TMP)<br>Precipitation (PRE)<br>Vapor (VAP)<br>Cloud Cover (CLD)<br>Wet Days (WET)<br>Frost Days (FRS) | Climate | CRU |
| Global Horizontal (SOL)<br>Direct Normal (SOL)<br>Global Extraterrestrial (SOL)<br>Direct Extraterrestrial (SOL) | Solar Radiation | NCDC |
| 1-year return level for temperature extreme (TMP.EXT) | Climate | Estimated using temp from CDIAC |



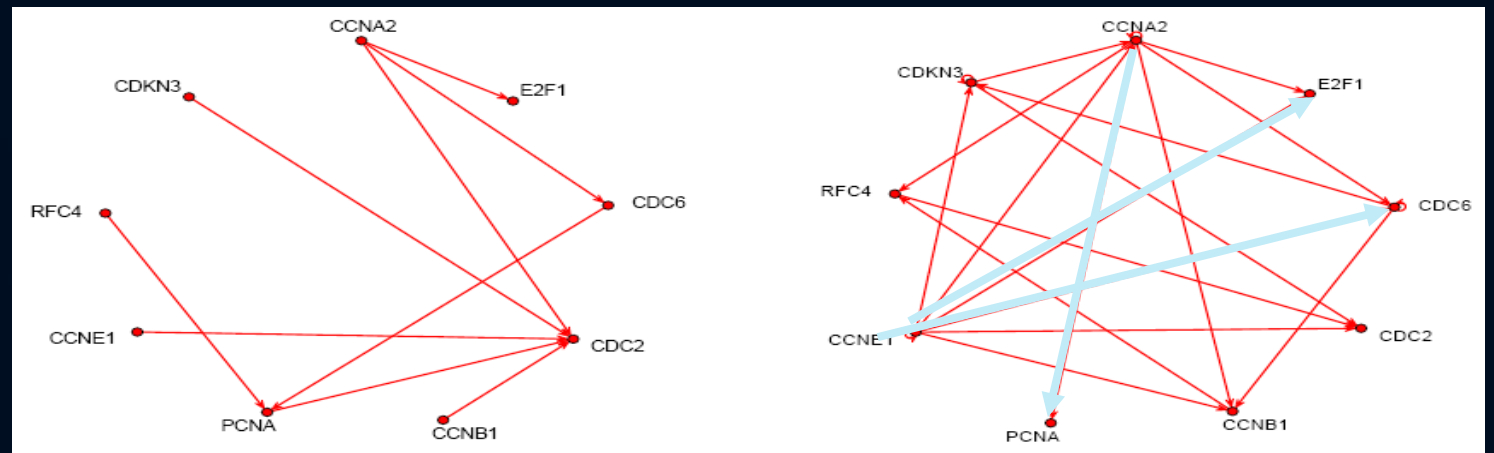*IPCC – Climate Change 2001, the scientific basis*

# Modeling Gene Regulations in Micro-array Data

- We applied our Grouped Graphical Granger Modeling technique to discover gene expression regulatory networks for the human cancer cell HeLa S3

- Our method achieved higher accuracy, and uncovered previously unknown relationships, given data involving 1134 genes

  - CCNA2 -> PCNA verified in [Liu, et al 2007]

  - CCNE1 -> ETF1 verified in [Merdzhanova, et al 2007]

  - CCNE1 -> CDC6 verified in [Furstenthal, et al 2001]



Full output network



Known links in database (BioGRID)        Causal links discovered by our method
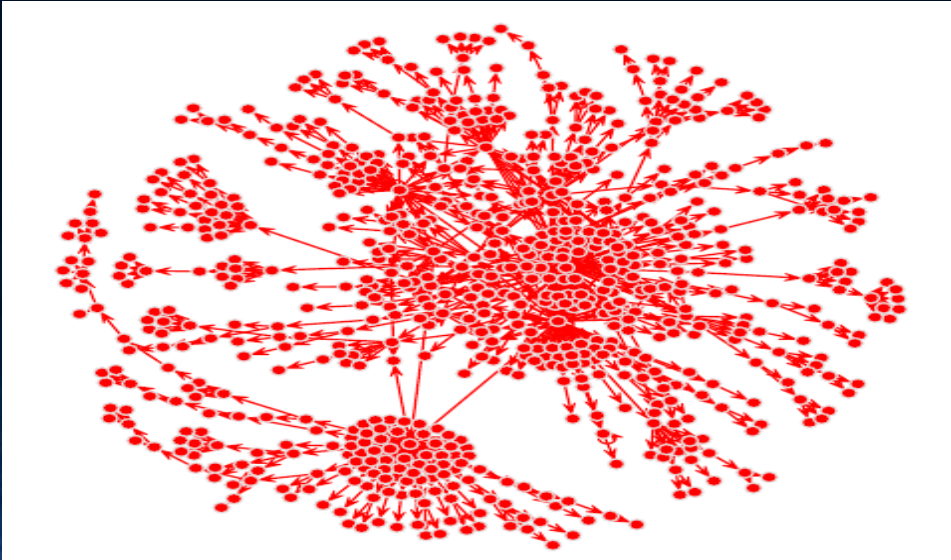
# Social Media Analytics: Key Influencer Identification
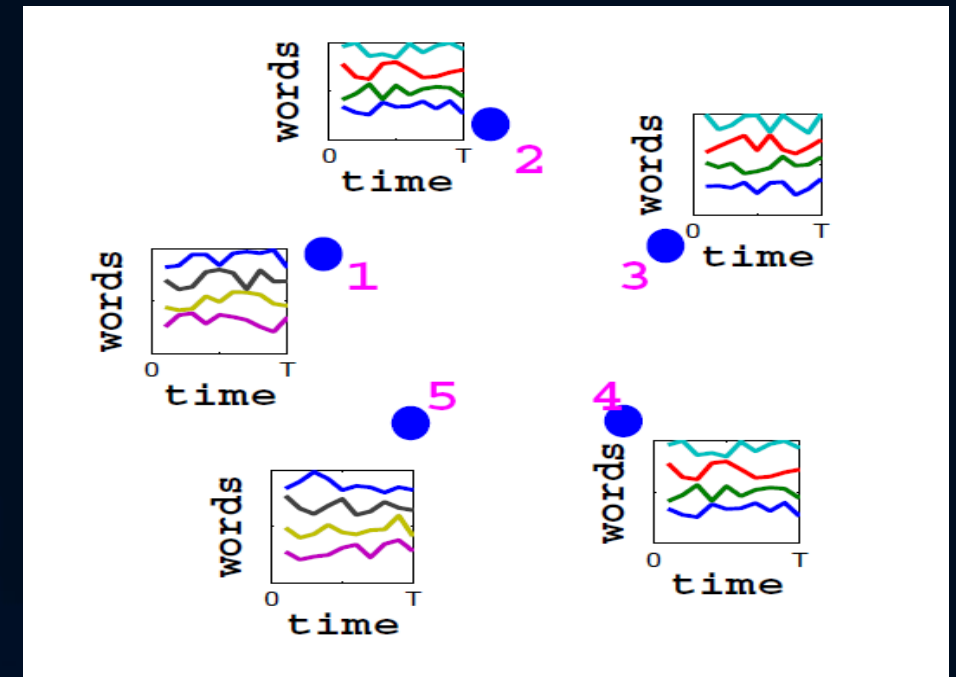
Question 1: Who is influencing whom?
Question 2: What is the context of influence?
Question 3: Who is most influential on a given topic?

Variables are Bloggers represented in terms of word usage across time





| Rank | GrangerPageRank |
|------|-----------------|
| 1 | http://www.theregister.co.uk |
| 2 | http://www.edbrill.com/ebrill/edbrill.nsf/ |
| 3 | http://soa.sys-con.com/ |
| 4 | http://lifehacker.com |
| 5 | http://www.dominoguru.com |
| 6 | http://communities.vmware.com |
| 7 | http://www.heise-online.co.uk |
| 8 | http://www.techworld.com |
| 9 | http://www.eweek.com |
| 10 | http://www.cmswire.com/ |

# Business KPI Monitoring and the CFO Performance Insights

- Captures significant, meaningful relationships between risk and performance metrics;

- Operationalizes the derived relationships to support various types of analysis

  - What-if scenario analysis

  - Goal seek

  - Root cause analysis

- Leverages multi-dimensional data for hierarchical modeling at various granularities

# IT Network Monitoring: Smart Cloud Predictive Insights

- Along with an "anomaly alarm", associated causal structure is presented, as well as "root cause" variables identified via alarm consolidation

# IBM Days (2000's – present)

- Smarter Government

  - Markov Decision Process and Tax Collections

- Smarter Planet , Enterprise & Cloud

  - Temporal Causal Modeling and Applications

- ***Smarter Agriculture***

  - ***Heterogeneous Data Analytics and Accelerated Plant Breeding***

# Smarter Agriculture

• Food production needs to be nearly doubled by 2050, to keep pace with population growth

• Challenged by limited resources, e.g. water and arable land, need to increase production efficiency



Food Gap



Source: : Rueters

Projected gap between demand and supply of row crops

# Department of Energy (DOE) ARPA-E Funded Project "Transportation Energy Resources from Renewable Agriculture (TERRA)"

- Project duration
  - 3 years (Sep. 2015 to Aug. 2018)

- Project Goals
  - To develop "**Automated Sorghum Phenotyping and Trait Development Platform**"
    - An automated high-throughput system for determining how variations in the sorghum genome impact field performance and agricultural productivity
    - Capacity to use sensing data from ground-based mobile and airborne platforms for automated phenotyping will advance plant breeding to maximize energy potential for transportation fuel

- Partnership
  - Purdue University, IBM Research, CSIRO



Gene expressions    SNP's    Drones    Phenomobiles

Genomic data    Field performance (phenotype) data

Trait development    Automated phenotyping

Genotype to Phenotype map

Plant breeding recommendation To maximize fuel energy potential

IBM

# Images from the Purdue Fields: Sorghum

IBM Research: Mobile, Solutions, and Mathematical Sciences

IBM

# Images from the Purdue Fields: Drones and Phenomobiles





IBM Research: Mobile, Solutions, and Mathematical Sciences

# Images from the Purdue Fields (videos)

- The IBM team (minus one) is here in the Purdue experimental field

- Here's the phenomobile in action
  - https://youtu.be/adCJDrB4Sus

# Project Scope

- **Overall Project: Complete Integrated Phenotyping Systems Solutions**

  - High Throughput Automated Hardware & Sensing Technologies (Purdue)
    - Optimize high-throughput remote-sensing technologies to acquire relevant data on sorghum plant phenotypes

  - Image/Sensor Data Analytics (Purdue & IBM)
    - Develop data analytics algorithms for segmentation and feature extraction

  - Genetics, Genomics and Bioinformatics (IBM)
    - Develop sophisticated genetic analysis pipelines to identify genes controlling sorghum performance

Overall deliverable will be "Automated Sorghum Phenotyping and Trait Development Platform"

IBM Research: Mobile, Solutions, and Mathematical Sciences

© 2015 IBM Corporation

# Image Data Feature Extraction by Machine Learning

- ## Manual Phenotyping (Purdue)

Direct measurements

Phenotypes/traits

- ## Automated Phenotyping (Purdue/IBM)

Plant height
Leaf number
Leaf width
Leaf length
Leaf area
Chlorophyll content
Fluorescence
Leaf, stem, panicle biomass
Stalk and root lodging
Harvest moisture
Biomass yield

Segmentation
Spectral feature extract

- ## Latent Phenotype Feature Learning (IBM)

Image data

Latent features

- Unsupervised Feature Learning Approaches

Indefinite phenotypes

?

- Multi-scale Modeling Approach

Specific phenotype

- Multi-task Feature Learning

Specific/multiple phenotypes

# Unified Genotype Phenotype Association Analysis

- Discovering underlying group structure in genotype phenotype association

SNPs          Phenotypes

**structured sparse models !**

- Discovering links between heterogeneous data (genotype, phenotype, etc)

Mutations (SNPs)    Gene expression    Phenotypic features

Plant height
Leaf number
Leaf width
Leaf length
Leaf area
Chlorophyll content
Fluorescence
Leaf, stem, panicle biomass
Stalk and root lodging
Harvest moisture
Biomass yield

Phenotypes

**Mixed graphical models !**

- Discovering indirect links with less reliable phenotype measurement data

**robust modeling methods !**

# Hyperspectral Data Analysis (Preliminary Results)

- Extracting hyperspectral signatures for forests, crops, soil, etc

**Example Data**
**Indian Pine Ground truth**

**Non-linear embedding & clustering**

**Average spectra**
**Per class**

**MI of K-means jumps up to 0.4128**

**Wavelengths**
**0.4 to 2.5 10^(-6) m**

IBM Research: Mobile, Solutions, and Mathematical Sciences

# Genotype Network & Phenotype Prediction (Preliminary Results)

Examples of selected SNPs and chromosome positions



Phenotype prediction accuracy improves by modeling multiple phenotypes

# Summary

Horizontal Scaling

Deployment

Prototype

Research

Research

Valley of Death

Valley of Death

Valley of Death

| NEC Research | Stochastic tree grammars | Protein structure prediction |
| RWC | On-line relation learning | Word clustering |
| | On-line active learning | Internet ad optimization | ADWIZ | NEC BU |
| | Active learning | Immunology | | NEC BU |

Gov Fund
Corp Res Fund
Business Fund

**NEC**

| FOAK | Targeted Marketing | Constrained MDP | CCOM | NYS TACOS | NBA Optimizer |
| | | FOAK | | Client | Sig |
| | IT monitoring | | | | |
| ER | TCM/STCM | Business KPI | CVAT for IBM | CFO Dashboard | Smart Cloud Predictive Insights |
| | FOAK | Social media analytics | JP | JP-Sig | JP |
| | | | | | SPSS Modeler |
| | | | | | JP |

| TERRA | GWAS | Image Analytics | Mixed graphical models | Robust models | Heterogeneous analytics |

IBM Research: Mobile, Solutions, and Mathematical Sciences

© 2015 IBM Corporation

IBM

# Some Take-Home Messages for Industrial Researchers

- ***Think long term, but be flexible !***

- **Decades:** Think of your career in phases (in the unit of decades)

- **Bottom-Up:** Long term strategic thinking is important, but they don't always come from top down strategy

  - Indulge yourself in basic research while you are young...

- **Breath-First:** Explore breath-first first and then go in depth

  - Focus on today's research and your path tomorrow will present itself

- **Challenge:** Do dare to cross the "valleys of death" sometime during your career

- **Adaptive funding:** Use a mixture of different funding types wisely

# Some Relevant Publications

**NEC days**

1. Predicting Protein Secondary Structure Using Stochastic Tree Grammars. Naoki Abe, Hiroshi Mamitsuka. *Machine Learning*, November 1997, Volume 29, Issue 2, pp 275-301

2. On-line Learning of Binary Lexical Relations Using Two-dimensional Weighted Majority Algorithms. Naoki Abe, Hang Li and Atsuyoshi Nakamura. Proceedings of The Twelfth International Conference on Machine Learning, July 1995.

3. The Lob-Pass Problem. Jun'ichi Takeuchi, Naoki Abe and Shun'ichi Amari. *Journal of Computer and System Sciences, 61(3), 2000.*

4. Learning to Optimally Schedule Internet Banner Advertisements. Naoki Abe and Atsuyoshi Nakamura. Proceedings of The Sixteenth International Conference on Machine Learning, July 1999.

5. Prediction of MHC Class I Binding Peptides by Dynamic Experiment Design based on Query Learning with Hidden Markov Models. Keiko Udaka, Hiroshi Mamitsuka, Yukinobu Nakaseko and Naoki Abe. *Journal of Immunology*, 169(10), 5744-5753, 2002.

# Some Relevant Publications (Cont'd)

Constrained Markov Decision Process

1. Optimizing debt collections using constrained reinforcement learning. Naoki Abe, Prem Melville, Cezar Pendus, Chandan K. Reddy, David L. Jensen, Vince P. Thomas, James J. Bennett, Gary F. Anderson, Brent R. Cooley, Melissa Kowalczyk, Mark Domick, Timothy Gardinier. KDD 2010: 75-84

2. Tax Collections Optimization for New York State. Gerard Miller, Melissa Weatherwax, Timothy Gardinier, Naoki Abe, Prem Melville, Cezar Pendus, David L. Jensen, Chandan K. Reddy, Vince P. Thomas, James J. Bennett, Gary F. Anderson, Brent R. Cooley. Interfaces 42(1): 74-84 (2012)

Temporal Causal Modeling

1. Spatial-temporal causal modeling for climate change attribution.  Aurelie C. Lozano, Hongfei Li, Alexandru Niculescu-Mizil, Yan Liu, Claudia Perlich, Jonathan R. M. Hosking, Naoki Abe. KDD 2009: 587-596

2. Grouped graphical Granger modeling methods for temporal causal modeling. Aurelie C Lozano, Naoki Abe, Yan Liu, Saharon Rosset, Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining, ACM, 2009.

3. Grouped graphical Granger modeling for gene expression regulatory networks discovery. Aurelie C Lozano, Naoki Abe, Yan Liu, Saharon Rosset, Bioinformatics, Oxford Univ Press, 2009.

TERRA Project

http://www.purdue.edu/newsroom/releases/2015/Q2/purdue-leading-research-using-advanced-technologies-to-better-grow-sorghum-as-biofuel.html