# Data-Driven Modeling and Optimization Algorithms for h-CPS

Hamsa Balakrishnan
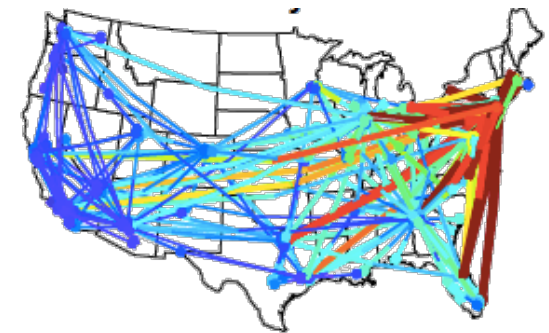
(with Jacob Avery, Michael Kasperski & Varun Ramanujam)

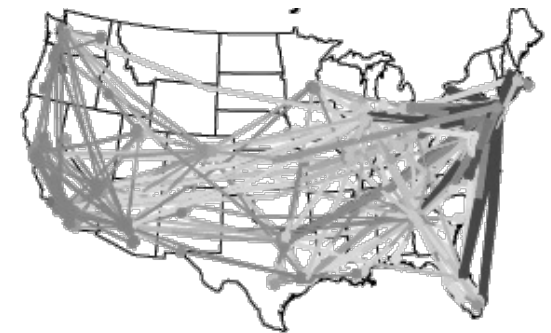Massachusetts Institute of Technology

# Most critical infrastructures are *h-CPS*

* Cyber + Physical + *human* (decision makers)
  * How do we model decision processes?
  * How do we estimate utility functions?
  * How do we predict system performance, when it depends on human decision-makers?
* Large-scale networks
  * Disruptions propagate through the system
* Multi-stakeholder systems
  * Optimization algorithms for resource allocation
  * Incentives for information-sharing



www.bls.gov





FORCES
FOUNDATIONS OF RESILIENT
CYBER-PHYSICAL SYSTEMS

6/11/14

# Most critical infrastructures are *h-CPS*

* Cyber **+** Physical **+** *human* (decision makers)
    * How do we model decision processes?
    * How do we estimate utility functions?
    * How do we predict system performance, when it depends on human decision-makers?
* Large-scale networks
    * Disruptions propagate through the system
* Multi-stakeholder systems
    * Optimization algorithms for resource allocation
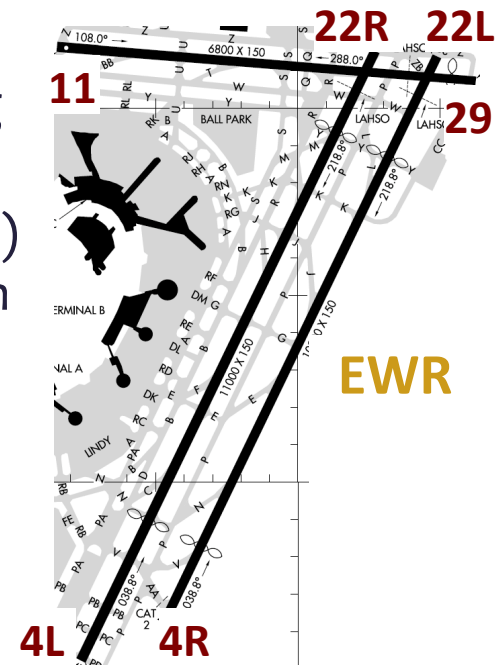    * Incentives for information-sharing

www.bls.gov

FORCES
FOUNDATIONS OF RESILIENT
CYBER-PHYSICAL SYSTEMS

6/11/14

# Modeling human decision processes

* Key challenge in h-CPS is modeling/predicting the behavior of the human participants
  * Discrete-choice models
    * Assuming decision-makers are rational, estimate their utility functions
    * Estimate relative weightings of different influencing factors
    * Use operational data (i.e., observations of decisions) to determine maximum-likelihood model of decision process
  * Approach demonstrated on airport configuration selection
    * Which runways are used for which operations
    * Primary driver of airport capacity



FORCES
FOUNDATIONS OF RESILIENT
CYBER-PHYSICAL SYSTEMS

# Discrete-choice models: Utility function

* Rich literature in transportation demand analysis (Ben-Akiva & Lerman 1985)
* Decision-maker chooses utility-maximizing option (from a set)
* Utility function is modeled as a linear function of the independent variables plus an error term

$$U_i = \underbrace{(\alpha_i + \beta_i \cdot X_i)}_{\text{Observed component, } V_i} + \underbrace{\epsilon_i}_{\text{Unobserved error}}$$

**Observed component, $V_i$**   **Unobserved error**

* For each observation, assume that the decision-maker chooses the alternative that maximizes utility, i.e., the choice $c_j$ such that

$$j = \underset{i:c_i \in C}{\arg\max}\, U_i$$

* Different models arise from assuming different forms of error term
  * Most widely used class of models assumes that the errors are independent and identically Gumbel distributed
  * Logistic Probability Unit, or Logit models
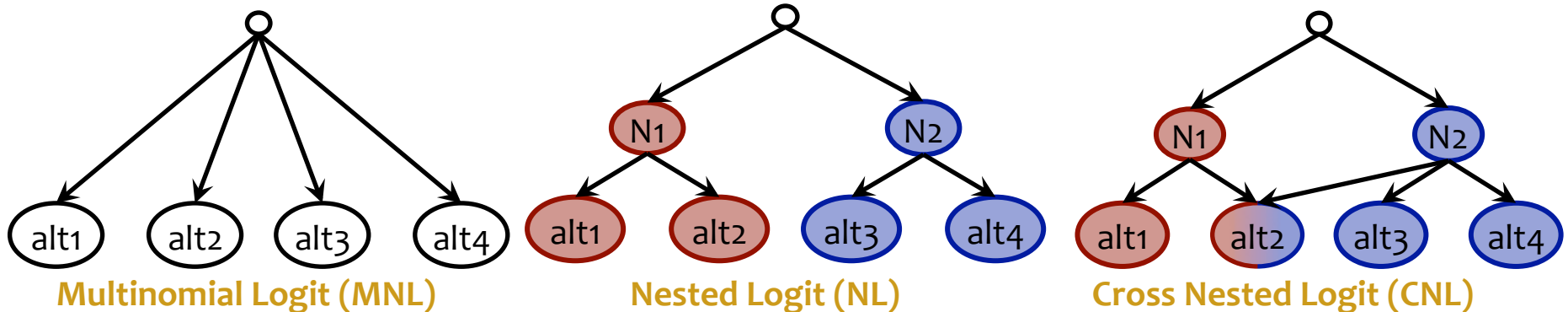
FORCES
FOUNDATIONS OF RESILIENT
CYBER-PHYSICAL SYSTEMS

# Discrete-choice models: Structure

* Alternatives (error terms) may not necessarily be independent

$$U_i = \underbrace{(\alpha_i + \beta_i \cdot X_i)}_{\text{Observed component, } V_i} + \underbrace{\epsilon_i}_{\text{Unobserved error}}$$

**Observed component, $V_i$**  **Unobserved error**

* Potential model structures:



**Multinomial Logit (MNL)**   **Nested Logit (NL)**   **Cross Nested Logit (CNL)**

* Alternatives within the same nest have correlated error terms

  * NL example: $V_{N1} = \dfrac{1}{\mu_{N1}} \log \sum_{j:c_j \in \{\text{alt1,alt2}\}} e^{\mu_{N1} V_j}; \ P(N1|\{N1, N2\}) = \dfrac{e^{V_{N1}}}{e^{V_{N1}} + e^{V_{N2}}}$

$$P(\text{alt1}|N1) = \frac{e^{\mu_{N1} V_{\text{alt1}}}}{\sum_{j:c_j \in \{\text{alt1,alt2}\}} e^{\mu_{N1} V_j}}$$

# Maximum-likelihood estimation

* Explanatory variables ($X_i$) of the utility function are determined iteratively
* For a given functional form of the utility function, the likelihood function of a given set of observations (over $N$ time periods, say) is

$$\mathscr{L}(\alpha,\beta) = P((c_1|C_1) \bigcap \ldots \bigcap (c_N|C_N)|\alpha,\beta,X)$$

  where $c_n$ is the choice observed at time $n$, and $C_n$ is the set of options
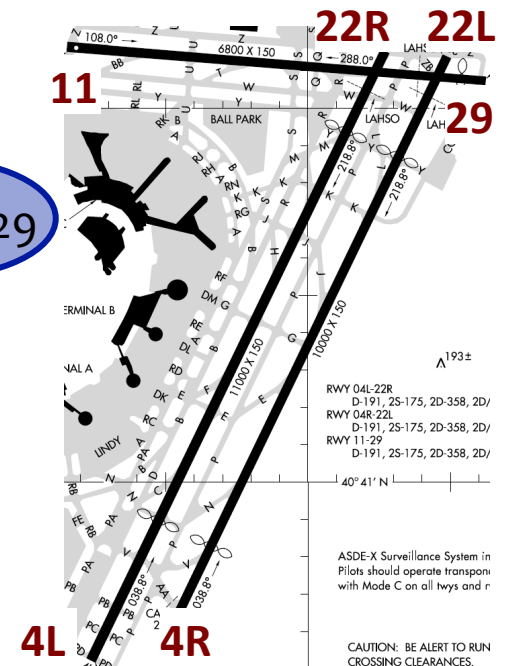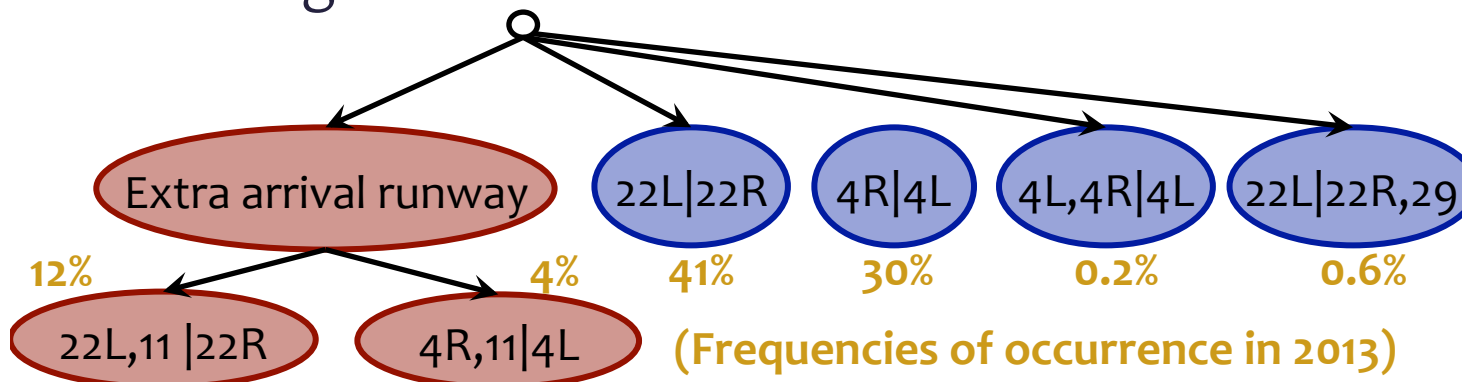
* Assuming that the observed choices at each time are conditionally independent given the explanatory variables, we get

$$\mathscr{L}(\alpha,\beta) = \prod_{i=1}^{N} P(c_i|C_i)$$

$$(\widehat{\alpha},\widehat{\beta}) = \arg\max_{\alpha,\beta} \mathscr{L}(\alpha,\beta)$$

* Nonlinear optimization problem (Bierlaire, 2003)
* Structure (MNL/NL/CNL) determined by checking statistical significance (Hausman-McFadden '84)

FORCES
FOUNDATIONS OF RESILIENT
CYBER-PHYSICAL SYSTEMS

6/11/14

# Results: Newark (EWR) airport case study

* Training data: 2011



* Extra arrival runway
* 22L|22R
* 4R|4L
* 4L,4R|4L
* 22L|22R,29

12%    22L,11 |22R    4%    4R,11|4L

41%    30%    0.2%    0.6%

(Frequencies of occurrence in 2013)

  * Only consider configurations used >50 times/year
* Validation data: 2013
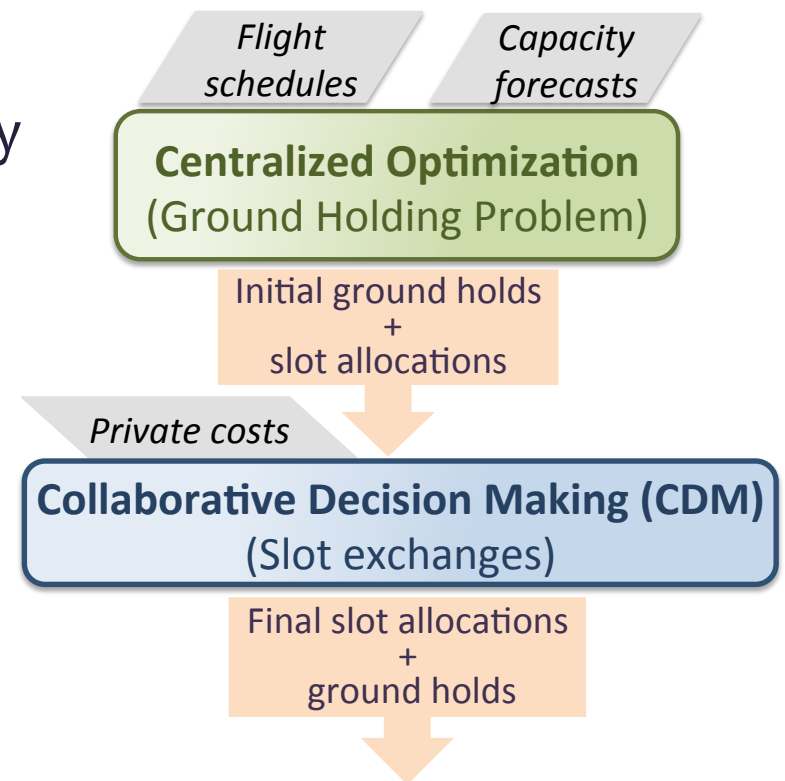  * Prediction accuracy: 73% (max possible: 89%)
  * 100% accuracy in distinguishing between use of the 22's and the 4's
* Currently extending approach to LGA and SFO
* Other applications

FORCES
FOUNDATIONS OF RESILIENT
CYBER-PHYSICAL SYSTEMS

6/11/14

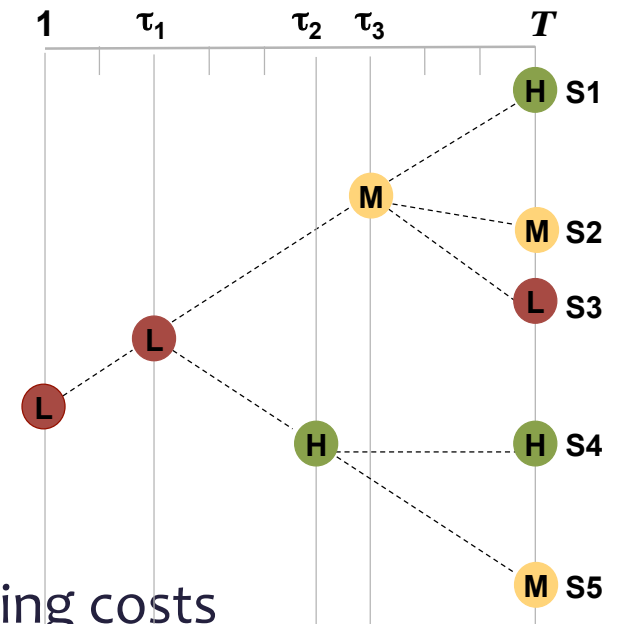# Resource Allocation: Optimization + Collaborative Decision Making

* Centralized optimization generally **assumes homogeneous delay costs**

* Airport capacity is uncertain, especially a few hours ahead of time

* Stochastic optimization formulations:

  * Static: Single-stage stochastic Integer Program (IP)

  * Dynamic: Multi-stage stochastic IP, differentiates between flights of different durations

  * *Hybrid:* Multi-stage stochastic IP, but does not differentiate between flights of different durations

Flight schedules     Capacity forecasts

**Centralized Optimization**
(Ground Holding Problem)

Initial ground holds
+
slot allocations

Private costs

**Collaborative Decision Making (CDM)**
(Slot exchanges)

Final slot allocations
+
ground holds

FORCES
FOUNDATIONS OF RESILIENT
CYBER-PHYSICAL SYSTEMS

# Static Ground Holding Problem

* Single-stage stochastic IP (Richetta & Odoni 1993)

$$\text{Minimize} \quad \sum_{\substack{n=0}}^{K} C_{g,n}\left(\sum_{t=1}^{T-n} A_{t,t+n}^{gq}\right) + \sum_{q \in Q} \pi_q\left(C_a \sum_{t=1}^{T} A_{q,t}^{aq}\right)$$

$$\text{subject to} \quad \sum_{j=t}^{t+K} A_{t,j}^{gq} = A_t^{dem}, \; \forall t \in \{1,..,T\}$$

$$A_{q,t}^{aq} \geq \sum_{j=t-K}^{t} A_{j,t}^{gq} + A_{q,t-1}^{aq} - A_{q,t}^{cap}, \; \forall t \in \{1,..,T\}, q \in Q$$

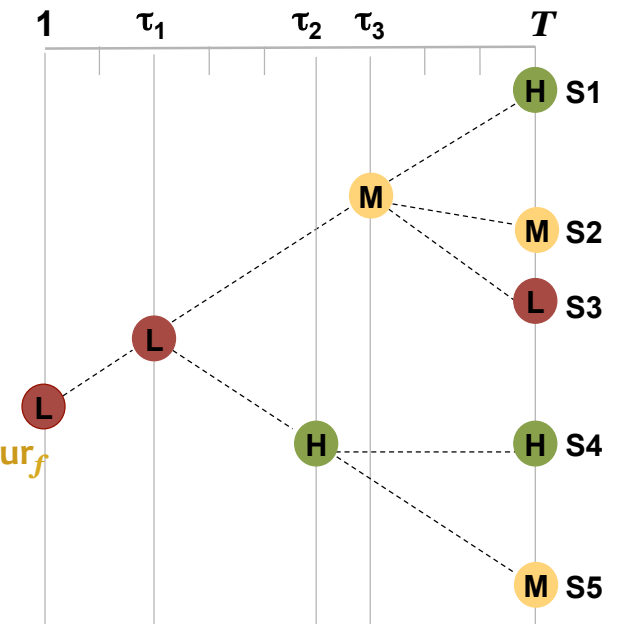$$A_{t,j}^{gq}, A_{q,t}^{aq} \in \mathbb{Z}^{+}, \; \forall t,j \in \{1,..,T\}, q \in Q$$



* LP relaxation is integer-optimal if ground-holding costs are marginally non-decreasing (Kotnyek & Richetta 2006)

FORCES
FOUNDATIONS OF RESILIENT
CYBER-PHYSICAL SYSTEMS

# Dynamic Ground Holding Problem

* Multi-stage stochastic IP (Mukherjee & Hansen 2007)

$$\text{Minimize} \sum_{q \in Q} \pi_q \left[ \sum_{f \in F} \left( \sum_{t=\text{arr}_f}^{\text{arr}_f+K} C_{g,t-\text{arr}_f} X_{f,t}^q \right) + \left( C_a \sum_{t=1}^{T} A_{q,t}^{\text{aq}} \right) \right]$$

subject to

$$\sum_{t=\text{arr}_f}^{\text{arr}_f+K} X_{f,t}^q = 1, \; \forall q \in Q, \forall f \in F$$

$$A_{q,t}^{\text{aq}} \geq \sum_{f \in F} X_{f,t}^q + A_{q,t-1}^{\text{aq}} - A_{q,t}^{\text{cap}}, \; \forall t \in \{1,..,T\}, q \in Q$$

$$X_{f,t}^{q_1} = X_{f,t}^{q_2}, \; \forall q_1, q_2 \in G_{t-\text{dur}_f} \quad \leftarrow \text{set of feasible scenarios at time } t\text{-dur}_f$$

$$X_{f,t}^q \in \{0,1\}, A_{q,t}^{\text{aq}} \in \mathbb{Z}^+, \; \forall t \in \{1,..,T\}, \forall q \in Q, \forall f \in F$$



* In general, LP relaxation solution is not integer-optimal

* $\mathscr{O}(FT^2 + T^2)$ integer decision variables

FORCES
FOUNDATIONS OF RESILIENT
CYBER-PHYSICAL SYSTEMS

6/11/14

# *Hybrid* Ground Holding Problem

* Multi-stage stochastic IP (Ramanujam & Balakrishnan CDC 2014, submitted)

$$\text{Minimize} \quad \sum_{\substack{q \in Q}} \pi_q \left( \sum_{n=0}^{K} C_{g,n} \sum_{t=1}^{T-n} X_{t,t+n}^q + C_a \sum_{t=1}^{T} A_{q,t}^{\text{aq}} \right)$$
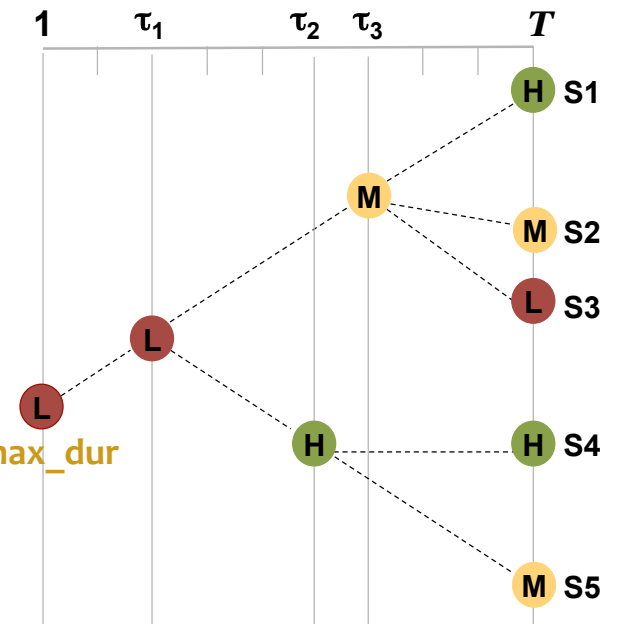
$$\text{subject to} \quad \sum_{j=t}^{t+K} X_{t,j}^q = A_t^{\text{dem}}, \; \forall t \in \{1,..,T\}, q \in Q$$

$$A_{q,t}^{\text{aq}} \geq \sum_{j=t-K}^{t} X_{j,t}^q + A_{q,t-1}^{\text{aq}} - A_{q,t}^{\text{cap}}, \; \forall t \in \{1,..,T\}, q \in Q$$

$$X_{t,j}^{q1} = X_{t,j}^{q2}, \; \forall q_1, q_2 \in G_{t-\text{max\_dur}}$$

$$X_{t,j}^q \in \mathbb{Z}^+, \; \forall t,j \in \{1,..,T\}, q \in Q$$

$$A_{q,t}^{\text{aq}} \in \mathbb{Z}^+, \; \forall t \in \{1,..,T\}, \forall q \in Q.$$

set of feasible scenarios at time $t$-max_dur

* In general, $\mathcal{O}(T^3)$ integer decision variables



FORCES
FOUNDATIONS OF RESILIENT
CYBER-PHYSICAL SYSTEMS

6/11/14

# Properties of the
# Hybrid Ground Holding Formulation

* Marginally non-decreasing ground-holding costs

LEMMA 1. *The hybrid stochastic SAGHP formulation yields an optimal solution with integer values for all variables $X_{a,b}^q$ ($\forall q \in Q; a, b \in \{1,..,T\}$) if the queue length variables ($A_{q,t}^{aq} \forall q \in Q, t \in \{1,..,T\}$) are constrained to have integer values, and the ground-holding costs are marginally non-decreasing (i.e., $C_{g,n+1} - C_{g,n} \geq C_{g,n} - C_{g,n-1} \forall n$).*

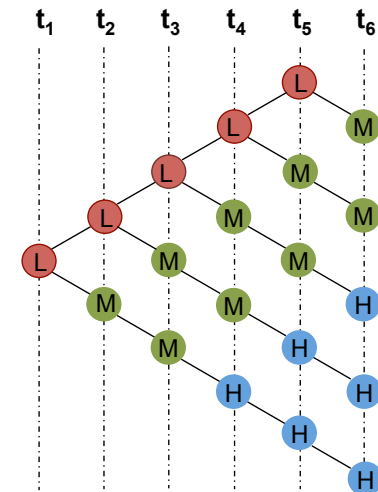* i.e., $\mathcal{O}(T^2)$ **integer variables instead of** $\mathcal{O}(T^3)$

[Ramanujam & Balakrishnan CDC 2014, submitted]

FORCES
FOUNDATIONS OF RESILIENT
CYBER-PHYSICAL SYSTEMS

6/11/14

# Properties of the
# Hybrid Ground Holding Formulation

* Marginally non-decreasing ground-holding costs + special scenario tree structure

LEMMA 2. *Given marginally non-decreasing ground-holding cost coefficients* $C_{g,n+1} - C_{g,n} \geq C_{g,n} - C_{g,n-1}$, $\forall n$, *and a capacity scenario tree forecast with sequentially non-decreasing capacity scenarios and sole element of uncertainty being time of improvement from lowest capacity state, the hybrid stochastic SAGHP formulation is guaranteed to have an integral optimum solution if the queue length variables for scenario* $T$ *(i.e.,* $A_{T,t}^{aq}$ $\forall t \in \{1,..,T\}$*) are constrained to be integers.*

* i.e., $\mathcal{O}(T)$ **integer variables instead of** $\mathcal{O}(T^3)$

[Ramanujam & Balakrishnan CDC 2014, submitted]



FORCES
FOUNDATIONS OF RESILIENT
CYBER-PHYSICAL SYSTEMS

6/11/14

# Optimization formulations of Collaborative Decision Making

## Intra-airline substitution

$$\text{Minimize} \quad \sum_{f_1 \in F_a} \sum_{f_2 \in F_a} C_{f_1,f_2} X_{f_1,f_2}$$

subject to:

$$\sum_{f_1 \in F_a} X_{f_1,f_2} = 1, \; \forall f_2 \in F_a$$

$$\sum_{f_2 \in F_a} X_{f_1,f_2} = 1, \; \forall f_1 \in F_a$$

$$X_{f_1,f_2} \leq \text{feas}_{f_1,f_2}, \; \forall f_1, f_2 \in F_a$$

$$X_{f_1,f_2} \in 0,1 \; \forall f_1, f_2 \in F_a$$

## Inter-airline substitution

$$\text{maximize} \quad \sum_{q \in Q} \pi_q \left[ \sum_{f \in F \backslash c} \sum_{r=1}^{T} \mathcal{B}(r) Y_{f,r}^q - M d_c^q \right]$$

subject to:

$$\sum_{t=\text{ETA}_f}^{\text{ETA}_f + K} X_{f,t}^q = 1, \; \forall q \in Q, \forall f \in F$$

$$A_{q,t}^{\text{aq}} \geq \sum_{f \in F} X_{f,t}^q + A_{q,t-1}^{\text{aq}} - A_{q,t}^{\text{cap}}, \; \forall t \in \{1,..,T\}, q \in Q$$

$$X_{f,t}^{q_1} = X_{f,t}^{q_2}, \; \forall q_1, q_2 \in G_{\text{ETA}_f - \text{max\_dur}},$$

$$d_c^q = \sum_{t=1}^{T} t X_{c,t}^q - k, \; \forall q \in Q,$$

$$\sum_{t=1}^{T} t X_{f,t}^q \leq \text{arr}_f^q, \; \forall q \in Q, f \in F \backslash c$$

$$A_{q,t}^{\text{aq}} \leq A q_{q,t}^{\text{aq,orig}}, \; \forall q \in Q, t \in \{1,..,T\}$$

$$Y_{f,r}^q = X_{f,\text{arr}_f^q - r+1}^q, \; \forall q \in Q, f \in F \backslash c, r \in \{1,..,T\}$$

$$X_{f,t}^q \in \{0,1\}, \; d_c^q \geq 0, \; \forall t \in \{1,..,T\}, \forall q \in Q, \forall f \in F$$

# Comparison of
# Ground Holding Problem Formulations

|  | Static | Hybrid | Dynamic |
|---|---|---|---|
| **Pre-CDM delay cost** | High (Worst) | Medium | Low (Best) |
| **Benefit from CDM** | High (Best) | Medium | Low (Worst) |
| **Equity** | High (Best) | Medium | Low (Worst) |
| **Tractability** | High (Best) | Medium | Low (Worst) |
| **Ease of implementation** | High (Best) | Medium | Low (Worst) |

[Ramanujam & Balakrishnan CDC 2014, submitted]

FORCES
FOUNDATIONS OF RESILIENT
CYBER-PHYSICAL SYSTEMS

# Summary

* Modeling human-driven decision processes is difficult, since utility functions are often not formally codified
  * Discrete-choice models present a way to determine utility functions as well as model structure
  * Data-driven models (descriptive, rather than prescriptive)
  * Maximum-likelihood estimation
* Multi-stakeholder optimization is a critical challenge in h-CPS
  * Different optimization formulations present tradeoffs in terms of
    * Computational tractability
    * System-optimal benefits
    * Incentives for participation
    * Incentives for information-sharing

FORCES
FOUNDATIONS OF RESILIENT
CYBER-PHYSICAL SYSTEMS