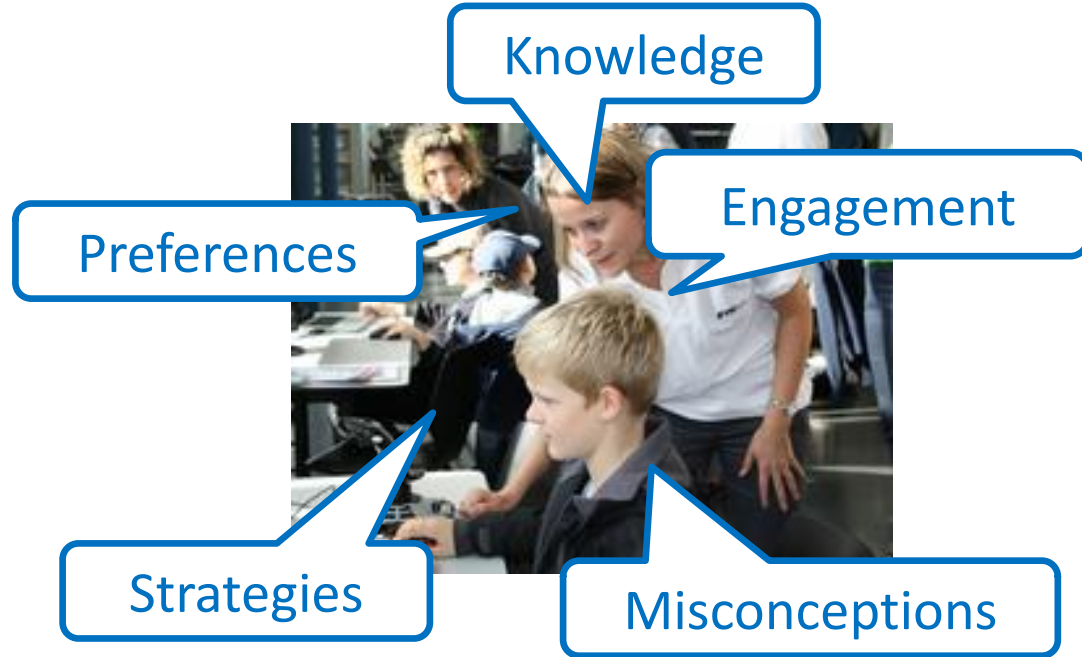# Modeling and Individualizing Learning in Computer-Based Environments

Tanja Käser

January 2020

SDSC

# Human teachers individualize learning

# Student models enable individualization



## Interaction
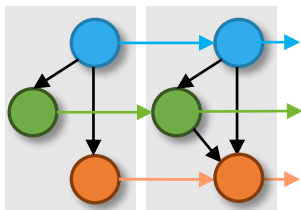
- Key stroke
- Mouse Click
- Speech
- Video

# Student models enable individualization



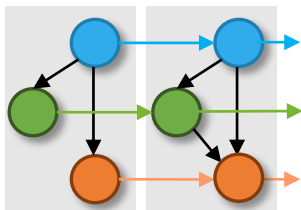**Interaction**

- Key stroke
- Mouse Click
- Speech
- Video

**Model**

- Detection
- Representation
- Prediction

# Student models enable individualization



**Interaction**

- Key stroke
- Mouse Click
- Speech
- Video

**Model**

- Detection
- Representation
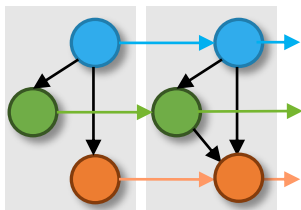- Prediction

**Individualization**

- Select new task
- Give feedback
- Provide hint

# Student models enable individualization



**Interaction**
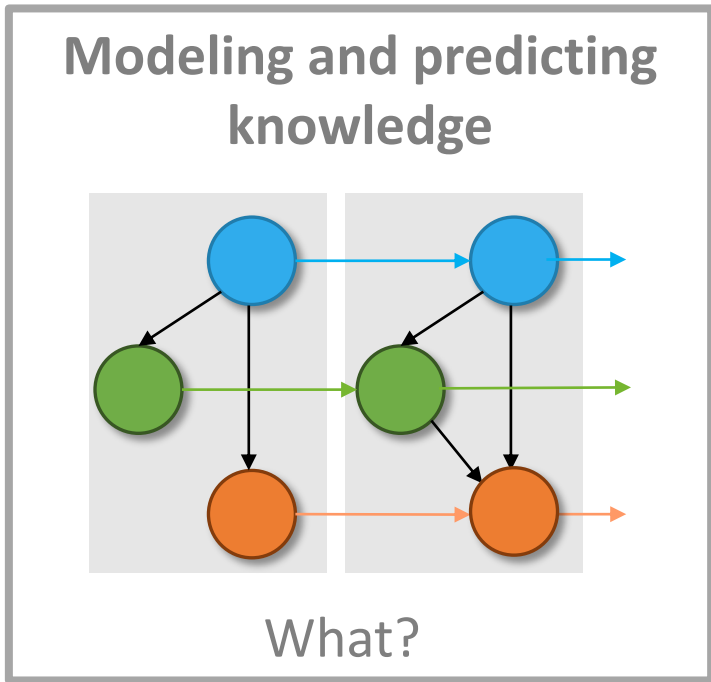
- Key stroke
- Mouse Click
- Speech
- Video

**Model**

- Detection
- Representation
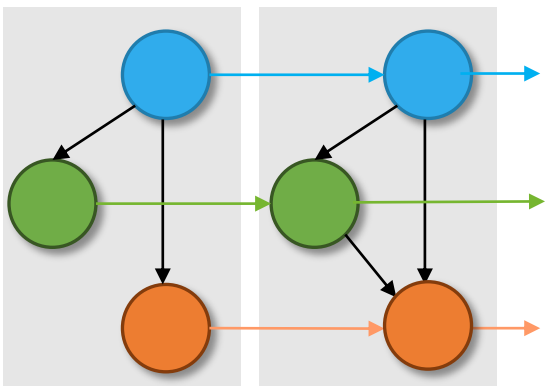- Prediction

**Individualization**

- Select new task
- Give feedback
- Provide hint

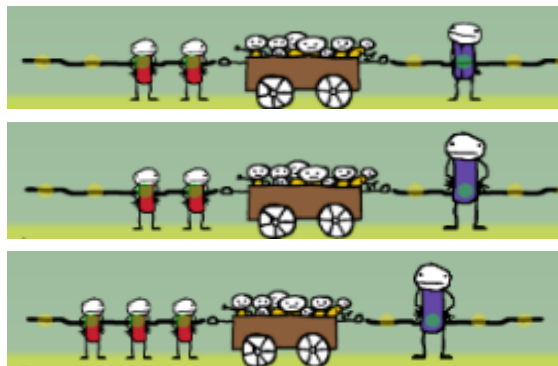# Modeling and Individualizing Learning in Computer-Based Environments

# Modeling and Individualizing Learning in Computer-Based Environments
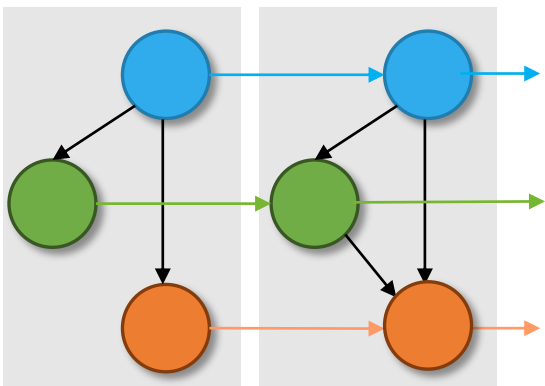


**Modeling and predicting knowledge**
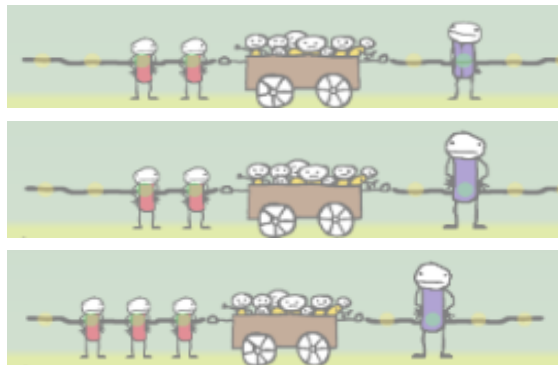
What?

**Detecting learner choices and strategies**

How?

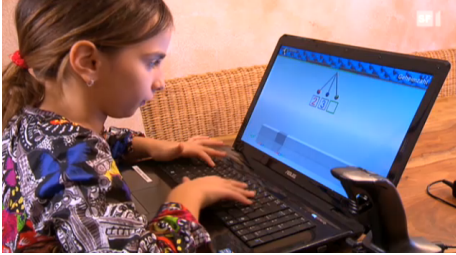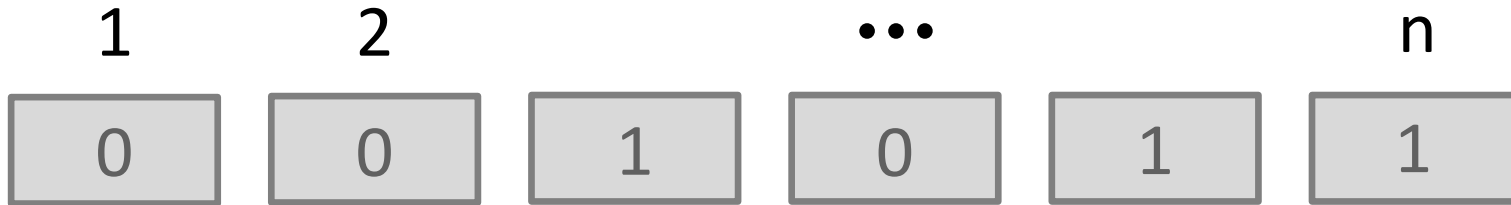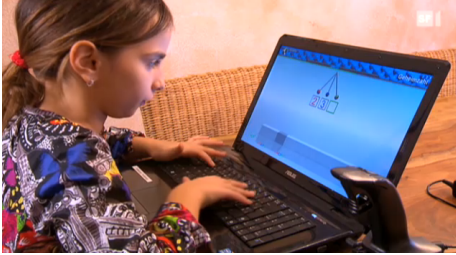# Modeling and Individualizing Learning in Computer-Based Environments
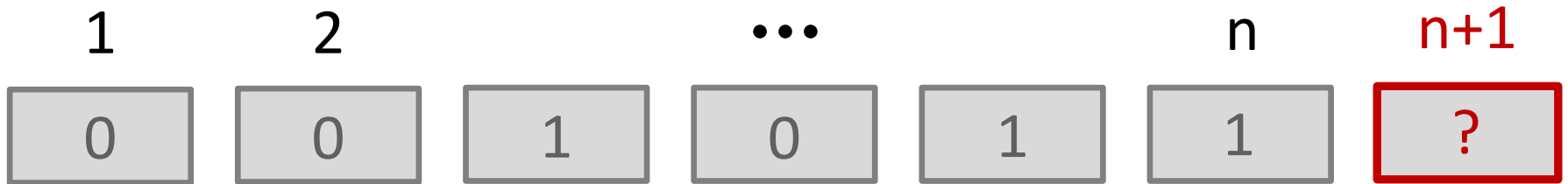
# Inferring knowledge based on student answers



**Subtraction 0-10**
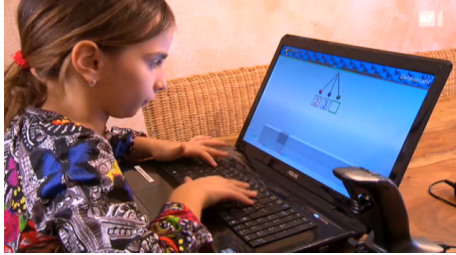
| 1 | 2 | ... | n |
|---|---|-----|---|
| 0 | 0 | 1 | 0 | 1 | 1 |

# Inferring knowledge based on student answers



**Subtraction 0-10**

| 1 | 2 | | $\cdots$ | | n | n+1 |
|---|---|---|---|---|---|---|
| 0 | 0 | 1 | 0 | 1 | 1 | ? |

# Bayesian Knowledge Tracing (BKT)



**Subtraction 0-10**

Latent variable

Observed variable

| 0 | 0 | 1 | 0 | 1 | 1 | ? |

# BKT models are simple, efficient, and interpretable

**Bayesian Knowledge Tracing (BKT)**



$$\theta = \{p_0, p_L, p_F, p_S, p_G\}$$

# BKT models are simple, efficient, and interpretable

**Bayesian Knowledge Tracing (BKT)**



Probability of learning or forgetting

$t = $     $t = 2$

$p_0 \rightarrow$ S1   $p_L, p_F$   S1

$p_S, p_G$

○ Latent variable

▢ Observed variable

$$\theta = \{p_0, p_L, p_F, p_S, p_G\}$$
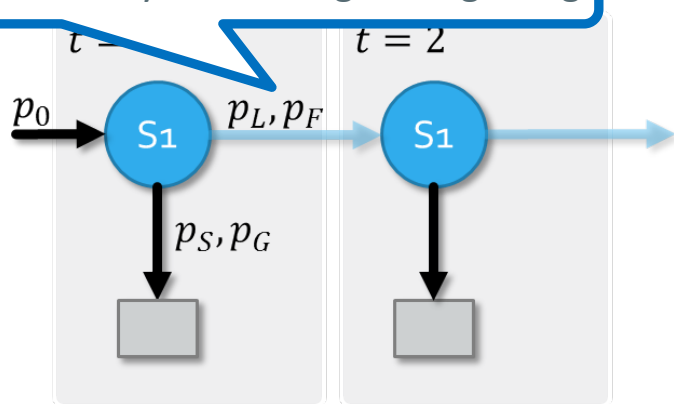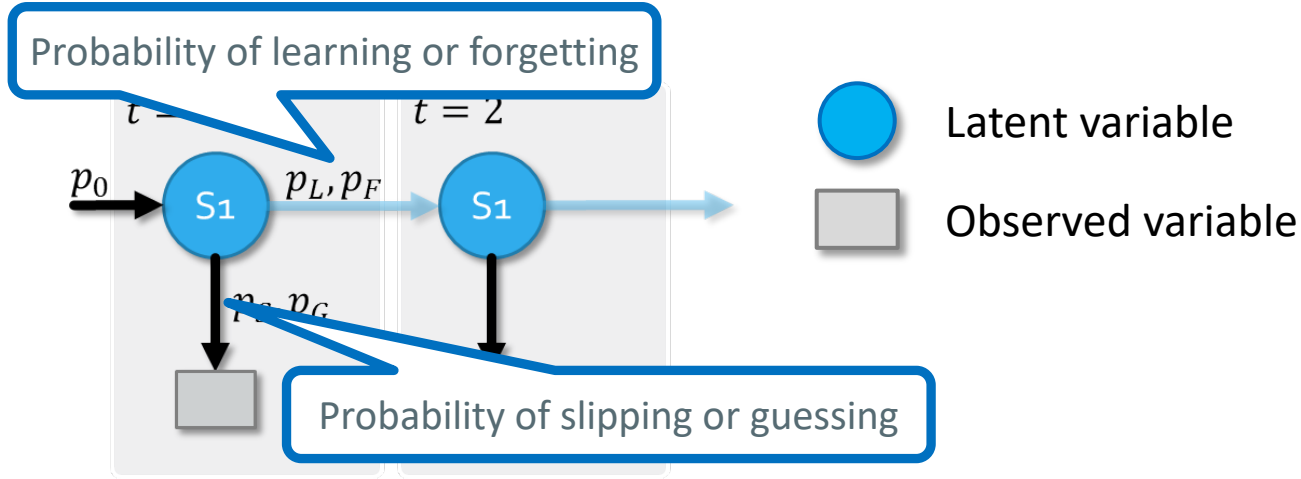
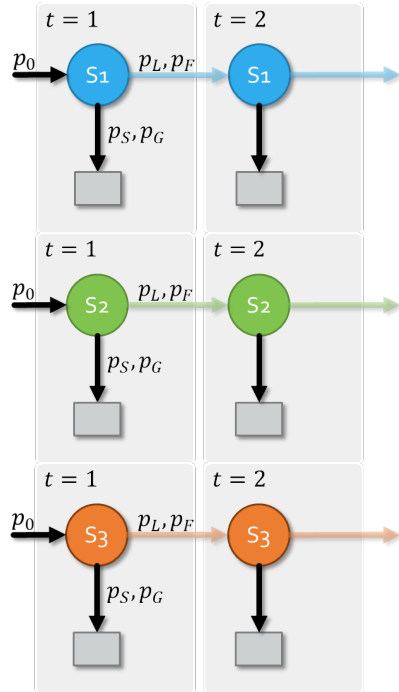# BKT models are simple, efficient, and interpretable

**Bayesian Knowledge Tracing (BKT)**



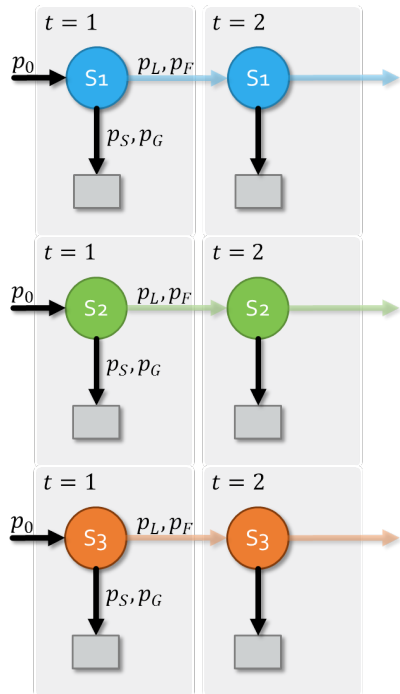$$\theta = \{p_0, p_L, p_F, p_S, p_G\}$$

# ... but they have limited representational power

**Bayesian Knowledge Tracing (BKT)**
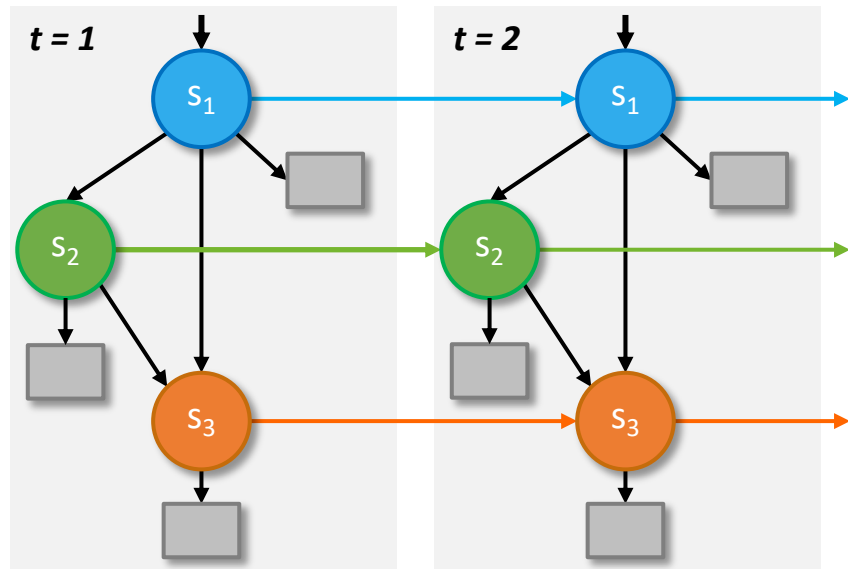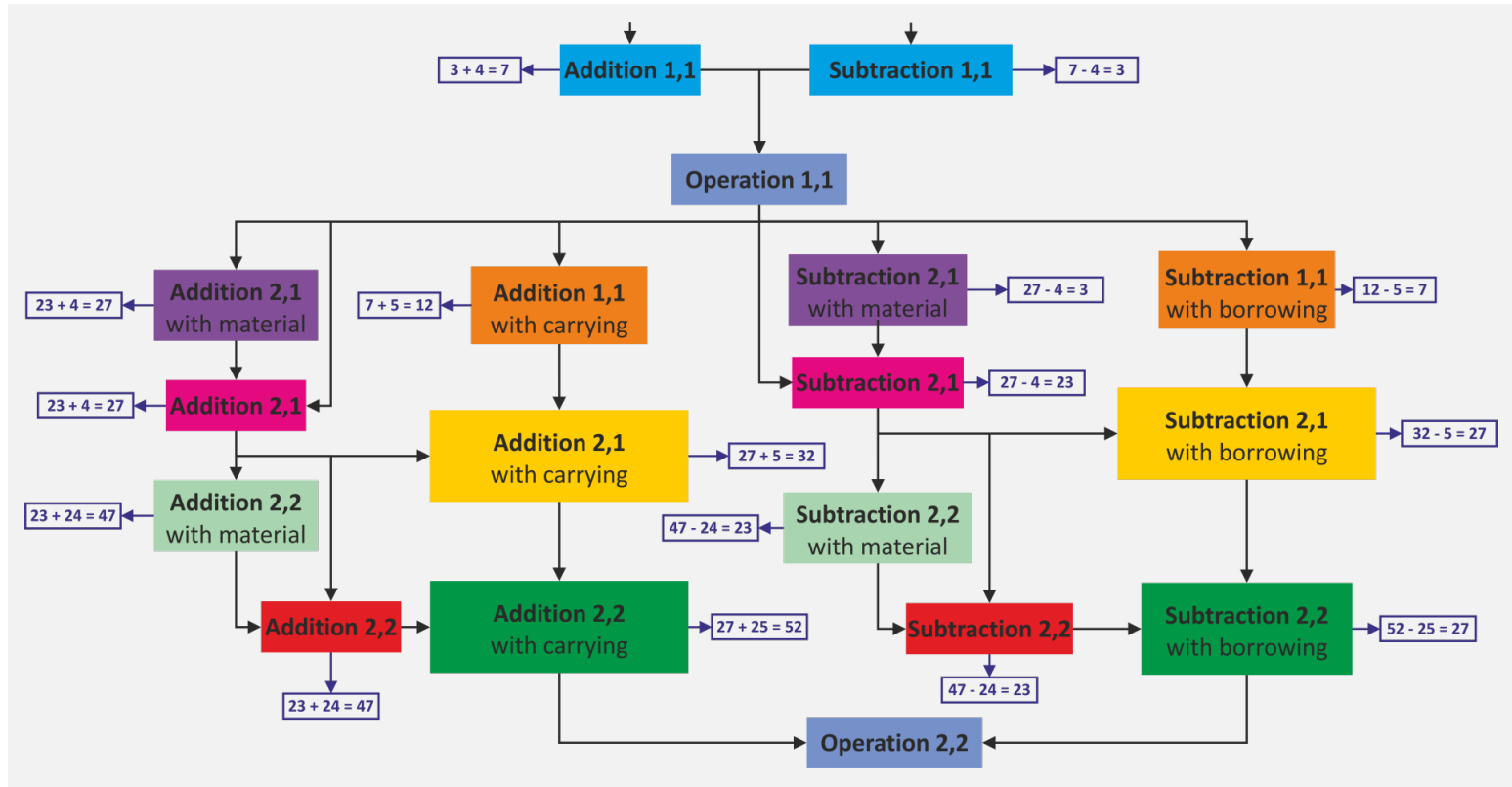
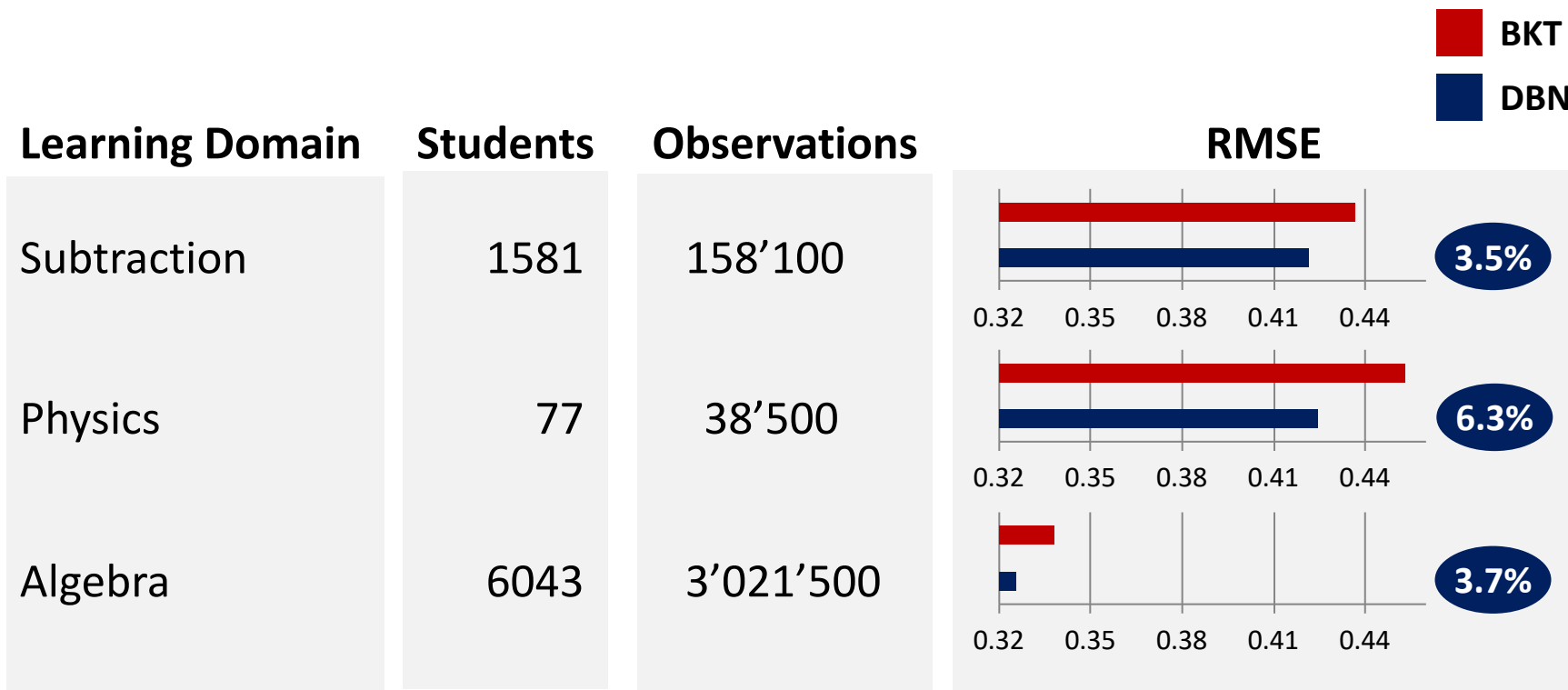# DBNs can model interactions between variables



**Bayesian Knowledge Tracing (BKT)**

**Dynamic Bayesian Networks (DBN)**

# Example: DBN representing mathematical skills



[Käser et al., Frontiers 2013; Käser et al., AISTATS 2014]

# DBNs outperform BKT in different learning domains



| Learning Domain | Students | Observations | RMSE |
|---|---|---|---|
| Subtraction | 1581 | 158'100 | 3.5% |
| Physics | 77 | 38'500 | 6.3% |
| Algebra | 6043 | 3'021'500 | 3.7% |

**BKT**

**DBN**

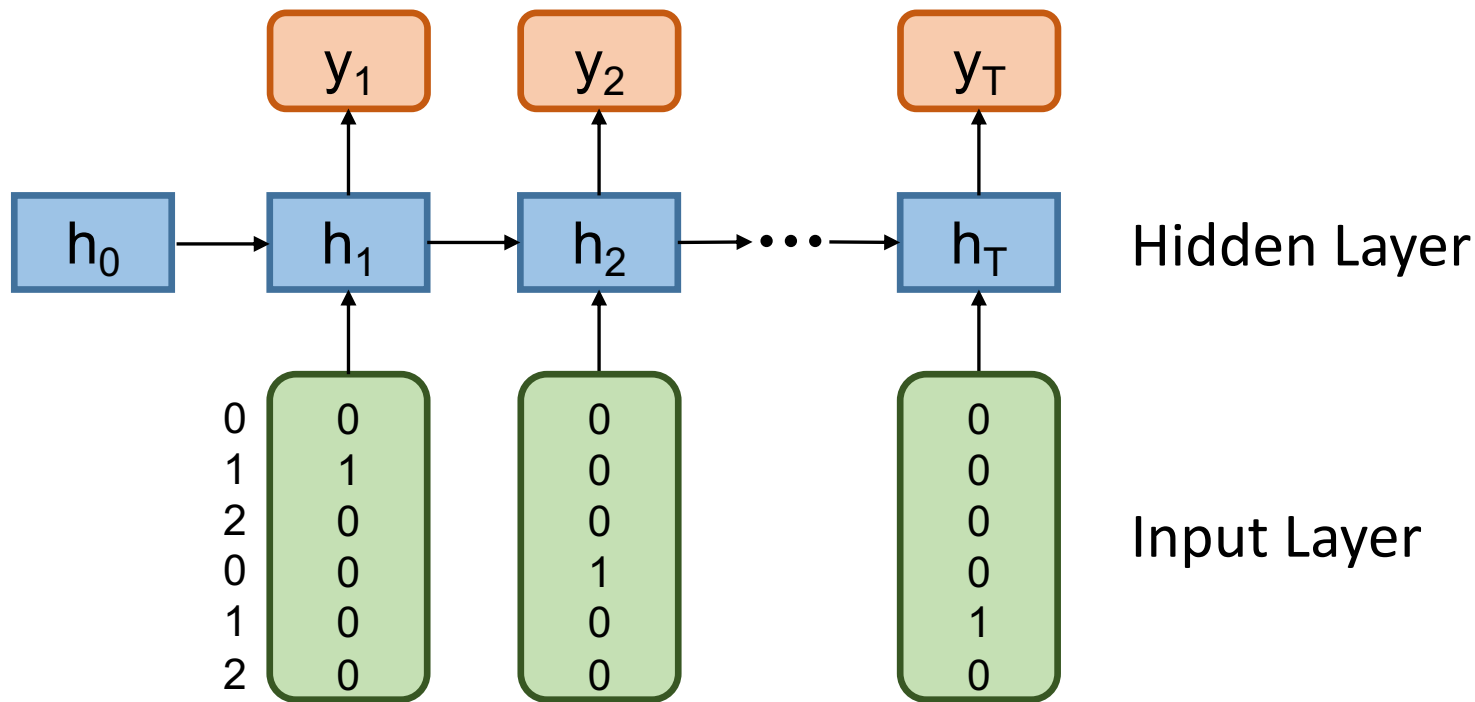[Käser et al., ITS 2014; Käser et al., IEEE TLT 2017]
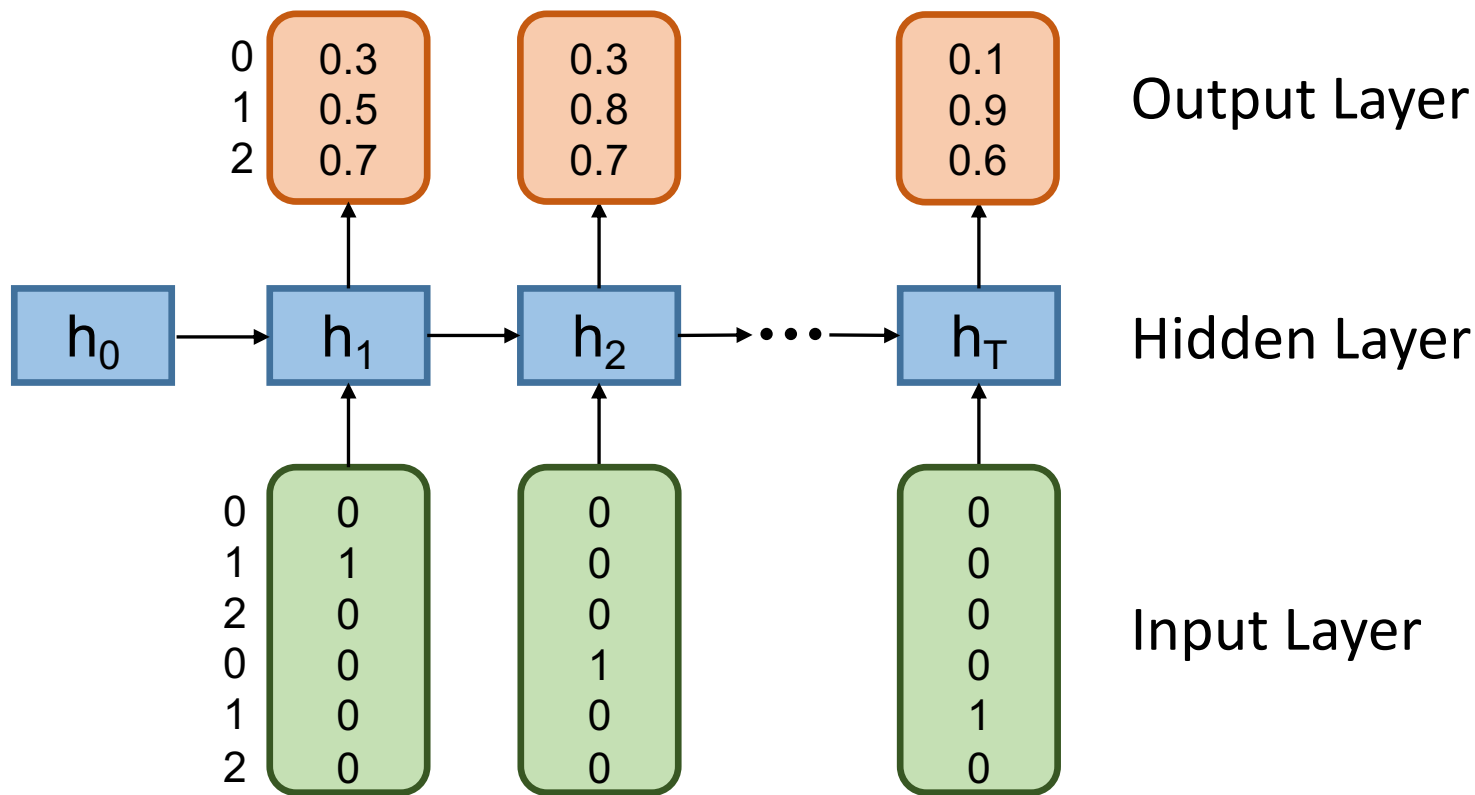
# Deep Knowledge Tracing

# Hidden layer captures relevant information



Hidden Layer

# Input layer represents observations
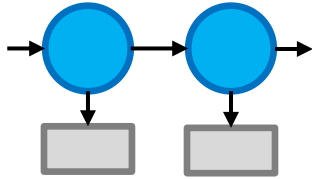
# Output layer consists of predicted probabilities



[Piech et al., NIPS 2015]

# Deep Knowledge Tracing outperforms BKT

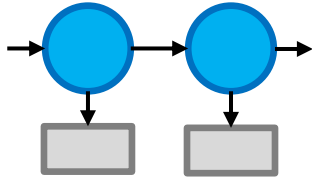| Data Set | Students | Observations | AUC |
|---|---|---|---|
| Khan Academy (Math) | 47'500 | 1'435'000 | |
| Assistments (Math) | 19'457 | 707'944 | |
| KDD Cup 2010 (Algebra) | 574 | 607'026 | |



BKT
DKT

[Piech et al., NIPS 2015; Xiong et al., EDM 2016]

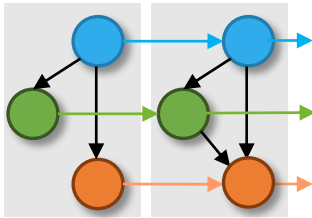# Modeling and Predicting Student Knowledge

Bayesian Knowledge Tracing is simple, efficient, and interpretable

# Modeling and Predicting Student Knowledge



Bayesian Knowledge Tracing is simple, efficient, and interpretable



Dynamic Bayesian Networks can represent the hierachical relations between the different skills
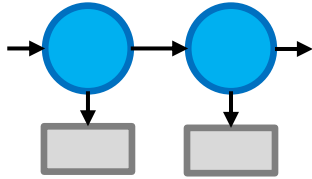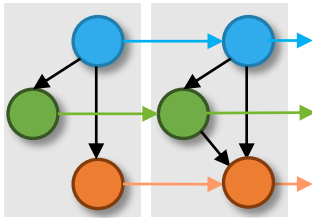
# **Modeling and Predicting Student Knowledge**

Bayesian Knowledge Tracing is simple, efficient, and interpretable

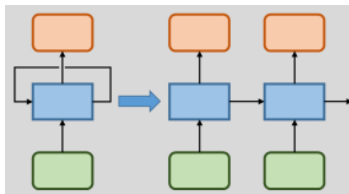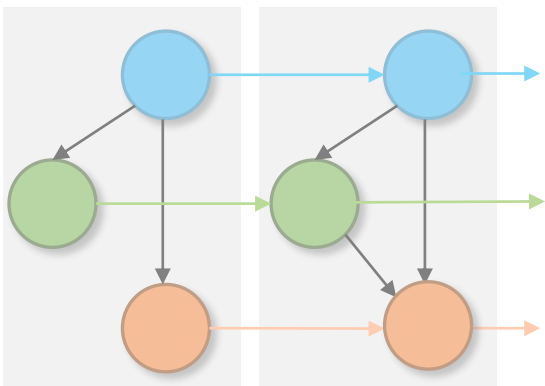Dynamic Bayesian Networks can represent the hierachical relations between the different skills

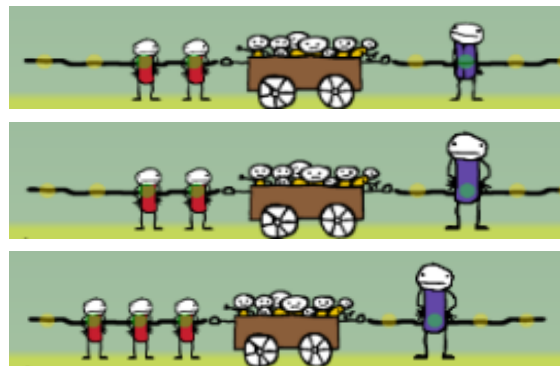Deep Knowledge Tracing can learn non-linear relationships and implicitly captures the relations between the skills

# Which team wins the tug-of-war?

# Students can freely choose between two modes

**Intro**

# Students can freely choose between two modes

**Intro**



**Challenge Mode**

Choose Winner

# Students can freely choose between two modes

**Intro**



**Challenge Mode**



Choose Winner

**wrong**



**Exploration Mode**

# Students can freely choose between two modes

**Intro**



**Challenge Mode**



Choose Winner

**wrong** | **correct**





**Exploration Mode**

# Students can freely choose between two modes

**Intro**



**Challenge Mode**

Choose Winner



**wrong**          **correct**

**Exploration Mode**



8 correct in a row

**Posttest**

# Students can be divided into six different clusters

**US School 1: 127 students**

# The best students explore systematically



US School 1: 127 students

# Persistent inquiry alone is not enough

**US School 1: 127 students**



Random Exploration

# Many students just try to beat the game



US School 1: 127 students

Trial-and-Error

# Adaptation based on students' learning behavior

# Exploring the use of recurrent neural networks

# LSTMs are similar or better at important levels



[Käser & Schwartz, EDM 2019]

# Questions?



tanja.kaeser@sdsc.ethz.ch

# References

1) Corbett, A. T., and Anderson, J. R. (1995*). Knowledge tracing: Modeling the acquisition of procedural knowledge*. User Modeling and User-Adapted Interaction

2) Yudelson, M.V., Koedinger, K.R., and Gordon, G.J. (2013). *Individualized bayesian knowledge tracing models*. Proceedings of AIED

3) Käser, T., Baschera, G., Kohn, J., Kucian, K. , Richtmann, V., Grond, U., Gross, M., and von Aster, M. (2013). *Design and evaluation of the computer-based training program Calcularis for enhancing numerical cognition.* Frontiers in Psychology

4) Käser, T., Klingler, S., Schwing, A., and Gross, M. (2014). *Computational Education using Latent Structured Prediction.* Proceedings of AISTATS

5) Käser, T., Klingler, S., Schwing, A., and Gross, M. (2014). Beyond KnowledgeTracing: Modeling Skill Topologies with Bayesian Networks. Proceedings of ITS

6) Piech, C., Bassen, J., Huang, J., Ganguli, S., Sahami, M., Guibas, L., and Sohl-Dickstein, J. (2015). *Deep Knowledge Tracing*. Proceedings of NIPS

7) Xiong, X., Zhao, S., Van Inwegen, E. G., Beck, J. E. (2016*). Going Deeper with Deep Knowledge Tracing*. Proceedings of EDM
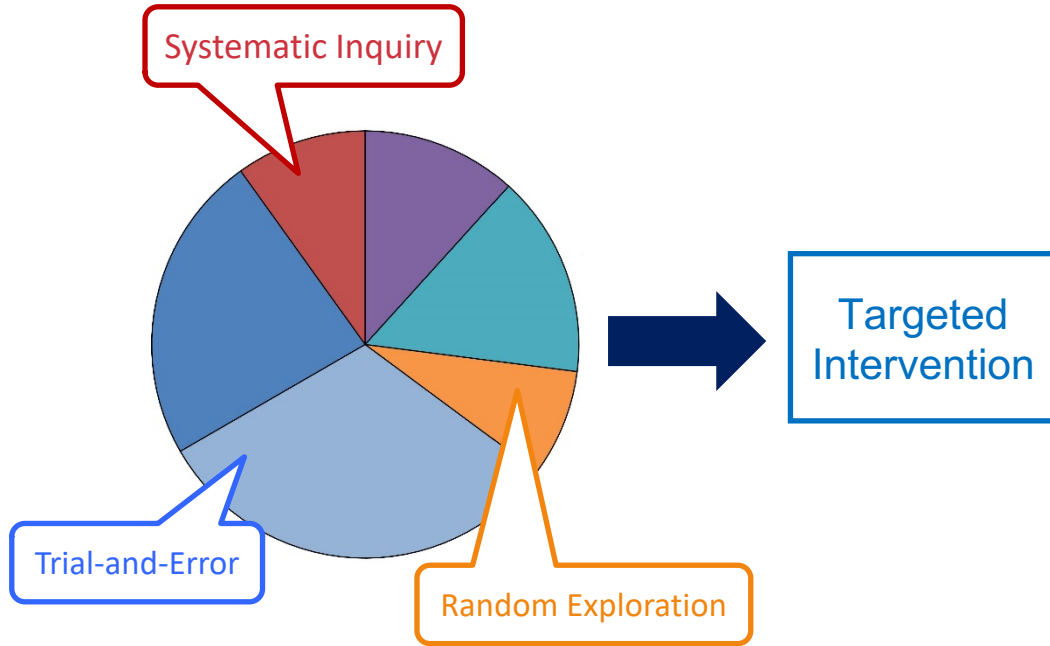
8) Käser, T., Klingler, S., Schwing, A., and Gross, M. (2017). Dynamic Bayesian Networks for Student Modeling*. IEEE Transactions on Learning Technologies*

9) Käser, T., and Schwartz, D. L. (2019). Exploring Neural Network Models for the Classification of Students in Highly Interactive Environments. *Proceedings of EDM*

# BACKUP

# Description of US data sets

|  | US School 1 | US School 2 |
|---|---|---|
| Number of students | 127 | 165 |
| Age | 8th grade | 8th grade |
| Time in exploration mode | 42% | 23% |
| Students passing the game | 87% | 97% |
| Students with perfect post-test | 24% | 34% |
| Average post-test score | 2.1 | 2.6 |

# Posttest

# Clustering students based on features describing their exploration behavior

➡️ **Number of challenge questions answered until passing a level (NC)**

# Clustering students based on features describing their exploration behavior

➡ Number of challenge questions answered until passing a level (NC)

➡ **Number of explored set-ups until passing a level (NS)**

# Clustering students based on features describing their exploration behavior

➡ Number of challenge questions answered until passing a level (NC)

➡ Number of explored set-ups until passing a level (NS)

➡ **Number of explored set-ups rated as strong until passing a level (NSS)**



Large = 3*Small

# The cluster solution was replicated on a second independent data set



US School 1: 127 students

US School 2: 165 students

Systematic Inquiry

Trial-and-Error

Random Exploration

Medium SES

High SES

# More students explore systematically



US School 1: 127 students

Systematic Inquiry

10%
12%
15%
8%
32%
23%

Medium SES

US School 2: 165 students

19%
9%
12%
9%
31%
21%

High SES

# Exploring students' inquiry strategies across cultural context

**US School 1: 127 students**



Medium SES

**US School 2: 165 students**



High SES

**Colombian Schools: 349 students**



Systematic Inquiry

Low-Medium SES

# Exploring students' inquiry strategies across cultural context

**US School 1: 127 students**



10%  12%  15%  8%  32%  23%

**Medium SES**

**US School 2: 165 students**



9%  12%  9%  31%  21%  19%

**High SES**

**Colombian Schools: 349 students**



3%  19%  5%  < 1%  60%  12%

**Trial-and-Error**

**Low-Medium SES**

# Clusters can be semantically interpreted



US School 1: 127 students

Medium SES

# Pairwise Clustering

Constant shift embedding transformation

similarities = distances in higher-dimensional Euclidean space

**k-Means Clustering**

# Computation of BIC

$$BIC = -2 \cdot \log(L) + k \cdot \log(n) + (k - 1) + 1$$

- L = likelihood of data
  - Fit Gaussian distribution per cluster
  - Estimate variance by distance to cluster centroid
  - Estimate mean by cluster centroid
  - Sum up gaussians over all clusters, taking into account the cluster probability
- k = number of clusters
- n = number of effective dimensions of transformation matrix

# Likelihood Computation

- Variance $\sigma^2 : \dfrac{1}{R-k} \cdot \sum_i (x_i - cc)^2$

  - R: Sample size

  - k: Number of clusters

  - cc: Centroid of according cluster

- $L_c = \dfrac{1}{p_c} \cdot \sum \dfrac{1}{\sqrt{2\pi\sigma^2}} \cdot e^{\left(\frac{x_i - cc}{\sigma}\right)^2}$

  - $p_c$: Prior probability for cluster

# Cluster Stability

- US School 1: Original data set

- US School 2: New data set

- Cluster US School 1 -> Original clustering solution (*OC*)

- k-Nearest Neighbor assigns each sample from school 2 to a cluster *c* of *OC* -> vector of predicted labels $l_p$

- Cluster US School 2 -> New clustering solution with labels $l_{NC}$

- Cluster stability = Hamming distance between $l_p$ and $l_{NC}$

# Exploring the use of recurrent neural networks

| | | 1 x 4 | 1 x 8 | 1 x 16 | 1 x 32 | 2 x 2 | 2 x 4 | 2 x 8 | 2 x 16 |
|---|---|---|---|---|---|---|---|---|---|
| LSTM | Predicting Sequence | | | | | | | | |
| GRU | Predicting Sequence | | | | | | | | |
| LSTM | Optimized for point in time | | | | | | | | |
| GRU | Optimized for point in time | | | | | | | | |

# Exploring the use of recurrent neural networks

|  |  | 1 x 4 | 1 x 8 | 1 x 16 | 1 x 32 | 2 x 2 | 2 x 4 | 2 x 8 | 2 x 16 |
|---|---|---|---|---|---|---|---|---|---|
| LSTM | Predicting Sequence |  |  |  |  |  |  |  |  |
| GRU |  |  |  |  |  |  |  |  |  |
| LSTM | Optimized for point in time |  |  |  |  |  |  |  |  |
| GRU |  |  |  |  |  |  |  |  |  |

# Output layer consist of predicted probabilities



Cluster Label

Output Layer

Hidden Layer(s)

Input Layer

$$y_t = \mathrm{softmax}(W_{yh}h_t + b_y)$$

# Model outputs a probability at each time step

# Model outputs a probability at the end

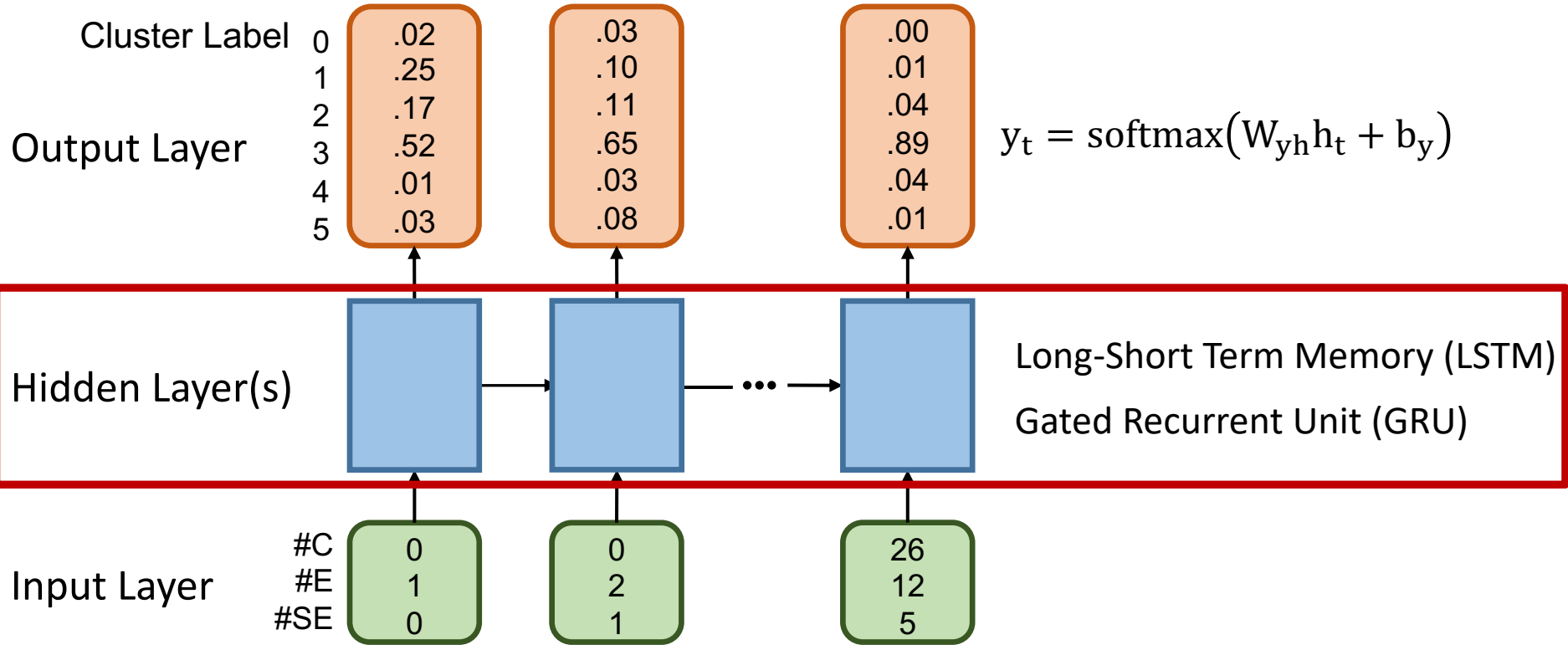# Hidden layer captures relevant information



Cluster Label

Output Layer

Hidden Layer(s)

Input Layer

$$y_t = \text{softmax}(W_{yh}h_t + b_y)$$

Long-Short Term Memory (LSTM)

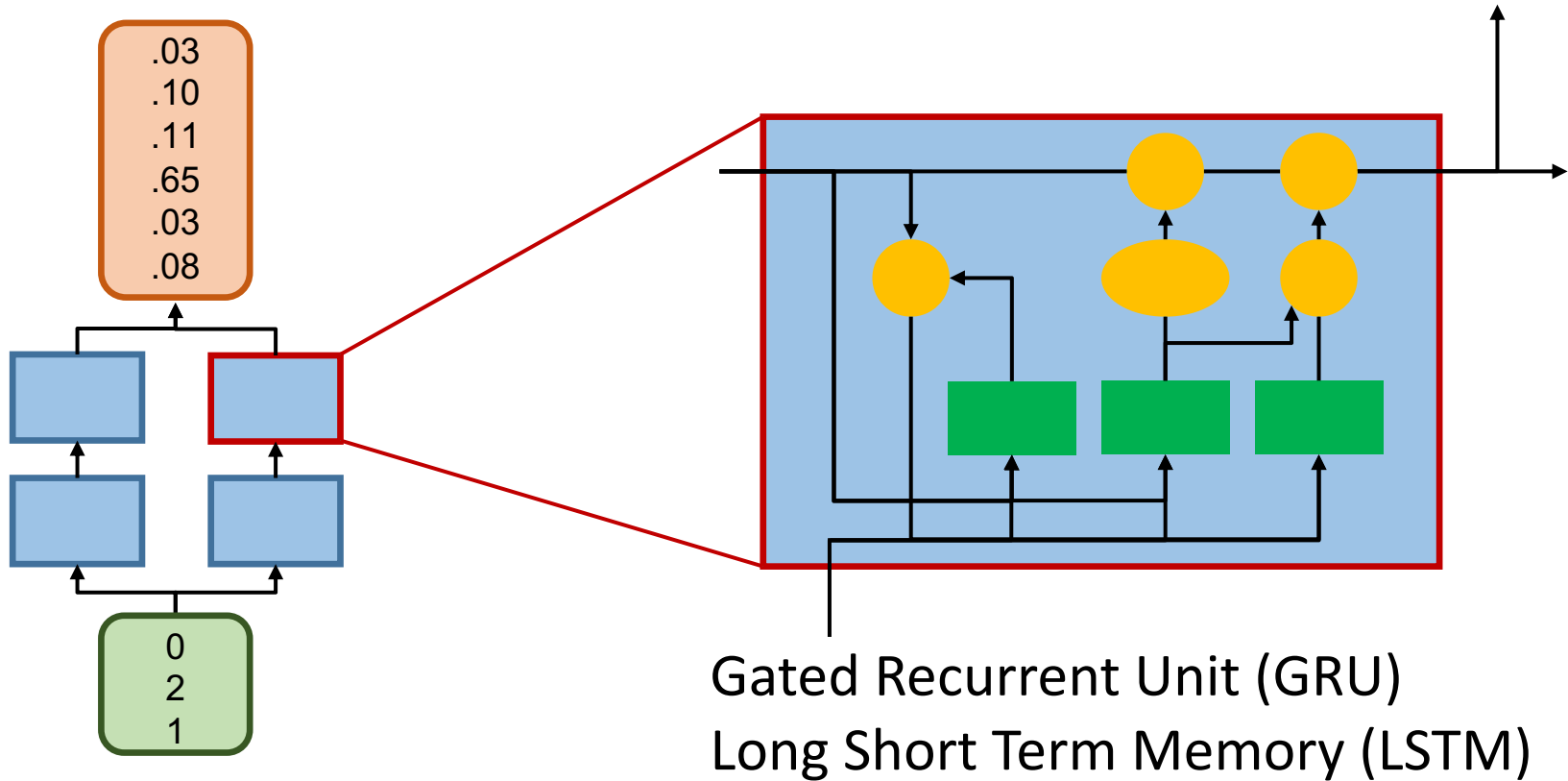Gated Recurrent Unit (GRU)

# Number of hidden layers and cells per layer vary

# Number of hidden layers and cells per layer vary

# Architecture of cells varies



Gated Recurrent Unit (GRU)
Long Short Term Memory (LSTM)

# Parameter learning is computationally intractable



$$\theta = \{p_1, p_2, ..., p_n\}$$

# Parameter learning is computationally intractable



$$\theta = \{p_1, p_2, \ldots, p_n\}$$

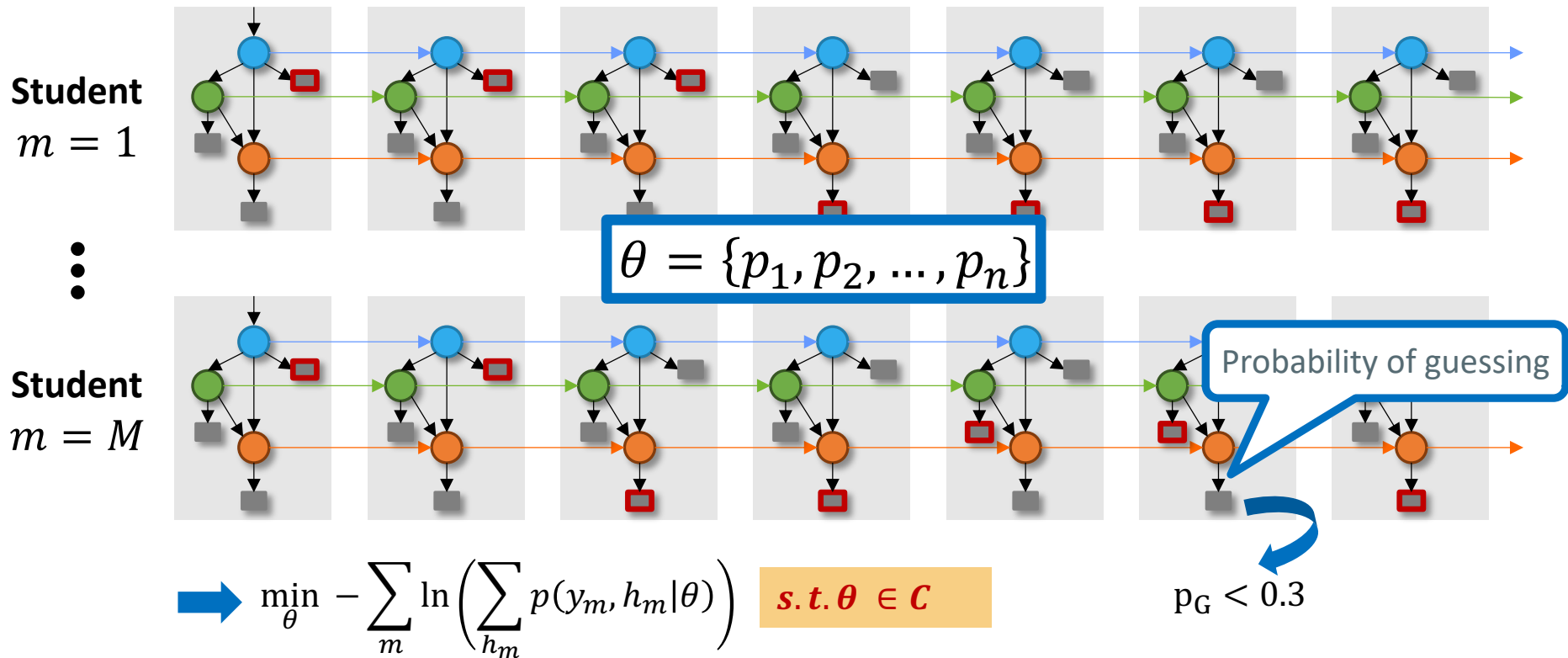$$\min_{\theta} - \sum_m \ln\left(\sum_{h_m} p(y_m, h_m | \theta)\right)$$

# Parameter constraints guarantee interpretability

# Parameter constraints guarantee interpretability



$$\theta = \{p_1, p_2, \ldots, p_n\}$$

Probability of guessing

Student $m = 1$

Student $m = M$

$$\min_\theta - \sum_m \ln\left(\sum_{h_m} p(y_m, h_m | \theta)\right) \quad s.t. \, \boldsymbol{\theta} \in \boldsymbol{C}$$

$p_G < 0.3$

**[Käser et al., AISTATS 2014]**

# From probabilistic notation to log-linear formulation

$$L(\theta) = \sum_m \ln \left( \sum_{h_m} p(y_m, h_m | \theta) \right)$$

$$L(w) = \sum_m \ln \left( \sum_{h_m} exp(\boldsymbol{w}^T \phi(y_m, h_{,m}) - \ln(Z)) \right)$$

[Käser et al., AISTATS 2014]
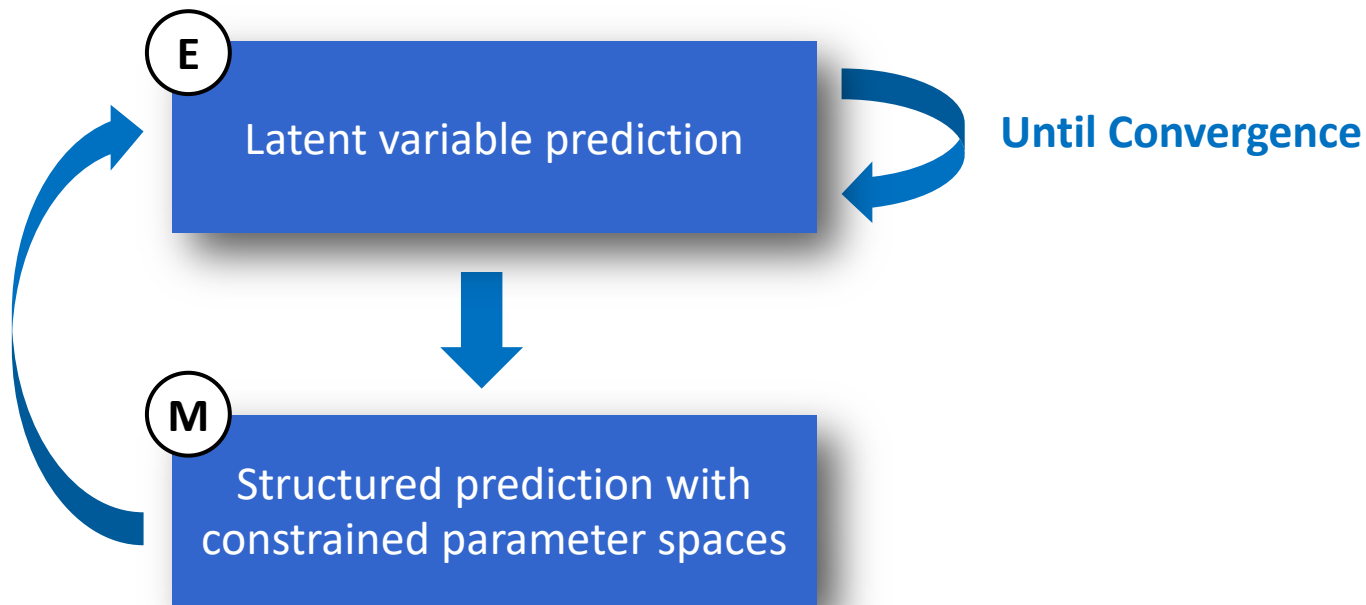
# From probabilistic notation to log-linear formulation

$$L(\theta) = \sum_m \ln \left( \sum_{h_m} p(y_m, h_m | \theta) \right)$$

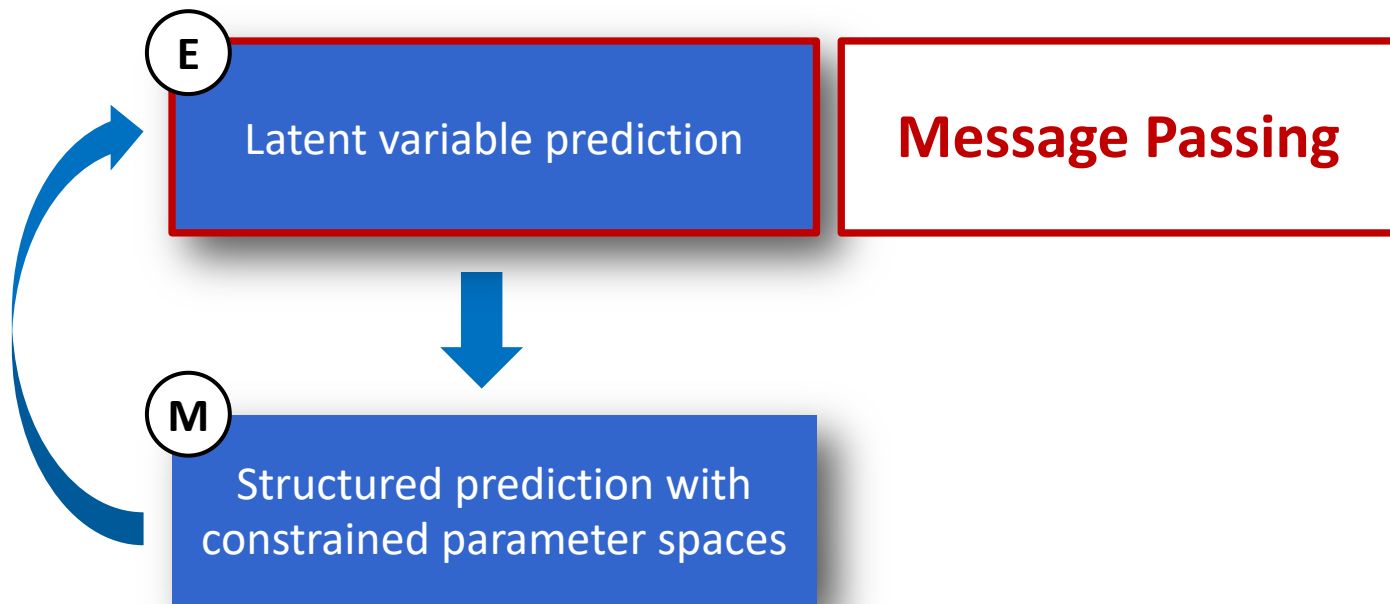$$\phi = 1 - 2v, V \in Y \cup H$$

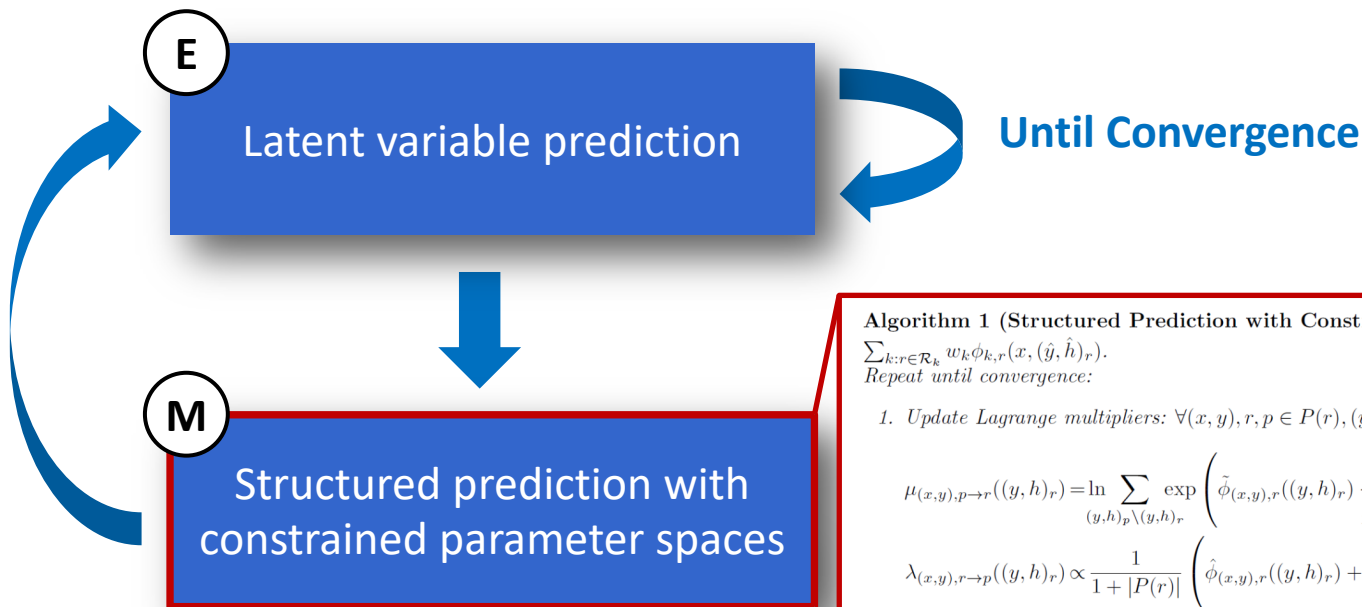$$L(w) = \sum_m ln \left( \sum_{h_m} exp(\mathbf{w}^T \phi(y_m, h_{,m}) - \ln(Z)) \right)$$

[Käser et al., AISTATS 2014]

# Constrained structured prediction with latent variables

# Constrained structured prediction with latent variables



[Schwing et al., ICML 2012]

# Constrained structured prediction with latent variables



[Käser et al., AISTATS 2014]

# Constrained structured prediction with latent variables



[Käser et al., AISTATS 2014]

# Constrained structured prediction with latent variables



[Käser et al., AISTATS 2014]

# DBNs outperform BKT in different learning domains

| Learning Domain | Students | Observations | RMSE | |
|---|---|---|---|---|
| | | |  | **BKT** **DBN** |
| Subtraction | 1581 | 158'100 | | 3.5% |
| Physics | 77 | 38'500 | | 6.3% |
| Algebra | 6043 | 3'021'500 | | 3.7% |
| Spelling | 7265 | 1'453'000 | | 0.7% |

[Käser et al., ITS 2014; Käser et al., IEEE TLT 2017]