Stochastic Process Model-Log Quality Dimensions an experimental study

Adam Burke, *Sander Leemans*, Moe Wynn, Wil van der Aalst and Arthur ter Hofstede



Event logs are stochastic
e.g. the log
[⟨a, b⟩<sup>20</sup>, ⟨a, b, c⟩<sup>2</sup>, ⟨a, b, c, c⟩<sup>1</sup>, ⟨e, f⟩<sup>1</sup>]
... already has frequency information



Simulation, analysis and recommendation need stochastic information

- Event logs are stochastic
- Control-flow models discard stochastic information
- Simulation, analysis and recommendation need stochastic information

- Event logs are stochastic
- Control-flow models discard stochastic information
- Stochastic process models retain stochastic information
- Simulation, analysis and recommendation need stochastic information







What dimensions describe the quality of stochastic process models?

Evaluation measures

 Earth-Movers' Stochastic Conformance

- Earth-Movers' Stochastic Conformance
- Entropy recall

- Earth-Movers' Stochastic Conformance
- Entropy recall
- Entropy precision

Exploration measures (13 new)

- Earth Movers' trace-wise (1)
- Probability mass (2)
- Fitness (6)
- Precision (2)
- Simplicity (3)
- Generalisation (4)

- Earth-Movers' Stochastic Conformance
- Entropy recall
- Entropy precision

Exploration measures (13 new)

- ► Earth Movers' trace-wise (1)
- Probability mass (2)
- Fitness (6)
- Precision (2)
- Simplicity (3)
- Generalisation (4)

Discover dimensions

- Earth-Movers' Stochastic Conformance
- Entropy recall
- Entropy precision

Exploration measures (13 new)

- Earth Movers' trace-wise (1)
- Probability mass (2)
- Fitness (6)
- Precision (2)
- Simplicity (3)
- Generalisation (4)

Discover dimensions

### Evaluation measures

- Earth-Movers' Stochastic Conformance
- Entropy recall
- Entropy precision

Identify dimensions

1. Use 6 public logs



- 1. Use 6 public logs
- 2. 9301 stochastic process models random, new genetic algorithm & discovered



- 1. Use 6 public logs
- 2. 9301 stochastic process models
- 3. 18 exploration measures



- 1. Use 6 public logs
- 2. 9301 stochastic process models
- 3. 18 exploration measures
- 4. Dimensional analysis



- 1. Use 6 public logs
- 2. 9301 stochastic process models
- 3. 18 exploration measures
- 4. Dimensional analysis



Dimensional Analysis 1: Correlations

Baselines: 2 log-only measures

### **Dimensional Analysis 1: Correlations**



Baselines: 2 log-only measures

## Dimensional Analysis 1: Correlations



- Baselines: 2 log-only measures
- Remove: 3 too-correlated measures: Trace Overlap Ratio, Trace Generalization Floor-1, Trace Generalization Floor-10

Find linear relation that best describes the data

- Find linear relation that best describes the data
- Find linear relation that best describes the data, orthogonal to first relation

- Find linear relation that best describes the data
- Find linear relation that best describes the data, orthogonal to first relation



- Find linear relation that best describes the data
- Find linear relation that best describes the data, orthogonal to first relation
- ▶ ... (15 times)



 15 linear combinations of measures

- 15 linear combinations of measures
- Scree plot: we choose 3



- 15 linear combinations of measures
- Scree plot: we choose 3



# Eigenvector Values For EX

Delete frame? See comment					
Standard deviations (1,, p=15):					
[1] 2.8598595882 1.6990123156 1.524660920	6 0.786566532	24 0.62972575	03 0.46003710	44 0.39432363	321 0.2941391
[14] 0.0149595210 0.0001853429					
Rotation (n x k) = (15 x 15):					
	PC1	PC2	PC3	PC4	PC5
ACTIVITY_RATIO_GOWER	-0.2313639	0.12995624	-0.28884572	-0.20701081	-0.83508956
TRACE_RATIO_GOWER_2	-0.2613284	0.21783178	-0.31657056	-0.17945506	0.09996934
TRACE_RATIO_GOWER_3	-0.25505	0.2061479	-0.34733325	-0.19427176	0.29392609
TRACE_RATIO_GOWER_4	-0.2508668	0.20311204	-0.33629168	-0.19388809	0.3374191
STRUCTURAL_SIMPLICITY_STOCHASTIC	-0.108313	0.45407014	0.32417104	-0.07154331	0.03656602
STRUCTURAL_SIMPLICITY_ENTITY_COUNT	-0.1220488	0.44260934	0.34036722	0.11738966	-0.03241646
STRUCTURAL_SIMPLICITY_EDGE_COUNT	-0.1306636	0.44564103	0.33960221	0.11777642	-0.0214105
TRACE_GENERALIZATION_DIFF_UNIQ	-0.2740367	-0.23386526	0.2618071	-0.31910582	0.05324514
EARTH_MOVERS_TRACEWISE	-0.2773889	-0.22840527	0.25235332	-0.3289332	0.05816939
TRACE_PROBMASS_OVERLAP	-0.2773891	-0.22840708	0.25235245	-0.32892646	0.05817131
ENTROPY_PRECISION_TRACEWISE	-0.3168787	0.04583473	-0.11641767	0.24698472	-0.02870561
ENTDODY EITNESS TDACEWISE	0.2000044	0.07931206	0 13210526	0 40540221	0.05756501

10

# Principal Components - Variation By Log



# Principal Components - By Model Generator



Remove random & genetic models



- Remove random & genetic models
- Add the 3 evaluation measures on EV models only



- Remove random & genetic models
- Add the 3 evaluation measures on EV models only



- Remove random & genetic models
- Add the 3 evaluation measures on EV models only
- Redo principal component analysis



#### Discovered dimensions



#### Identified dimensions



Play-out Entropy Project Fitness (HJFT) Generalization by Trace Floor 5 (TGE5) Play-out Entropy Project Precision (HJPT) Play-out Entropy Intersection Fitness (HIFT) Play-out Entropy Intersection Precision (HIPT) Trace Probability Mass Overlap (TMO) Earth Movers With Play-out Trace (EMT) Generalization by Trace Uniqueness (TGDU) Structural Simplicity by Edge Count (SSEDC) Structural Simplicity by Entity Count (SSENC) Structural Simplicity incl. Stochastic (SSS) Trace Ratio Gower 4 (TRG4) Trace Ratio Gower 3 (TRG3) Trace Ratio Course 2 (TRC2) Activity Ratio Gower (ARG) Earth Movers Truncated (tEMSC0.8) Entropy Precision (H P) Entropy Recall (H\_F) Trace Overlap Ratio (TOR) Generalization by Trace Floor 1 (TGF1) Generalization by Trace Floor 10 (TGF10)



Scaled contributions per component 🗆 min min min max

#### Discovered dimensions



#### Identified dimensions







#### Discovered dimensions



#### Identified dimensions





14

i min mid max

#### Discovered dimensions



#### Identified dimensions



ENTROPY



## Three Empirical Dimensions

Adhesion

How little effort is required to transform one stochastic language into another

Play-out Entropy Project Fitness (HJFT) Generalization by Trace Floor 5 (TGF5) Play-out Entropy Project Precision (HJPT) Play-out Entropy Intersection Fitness (HIFT) Play-out Entropy Intersection Precision (HIPT) Trace Probability Mass Overlap (TMO) Earth Movers With Play-out Trace (EMT) Generalization by Trace Uniqueness (TGDU) Structural Simplicity by Edge Count (SSEDC) Structural Simplicity by Entity Count (SSENC) Structural Simplicity incl. Stochastic (SSS) Trace Ratio Gower 4 (TRG4) Trace Ratio Gower 3 (TRG3) Trace Ratio Gower 2 (TRG2) Activity Ratio Gower (ARG) Earth Movers Truncated (tEMSC0.8) Entropy Precision (H P) Entropy Recall (H F) Trace Overlap Ratio (TOR) Generalization by Trace Floor 1 (TGF1) Generalization by Trace Floor 10 (TGF10)





## Three Empirical Dimensions

### Adhesion

How little effort is required to transform one stochastic language into another

### Entropy

The amount of information in a system

In this case, the combination of log and model

Play-out Entropy Project Fitness (HJFT) Generalization by Trace Floor 5 (TGF5) Play-out Entropy Project Precision (HJPT) Play-out Entropy Intersection Fitness (HIFT) Play-out Entropy Intersection Precision (HIPT) Trace Probability Mass Overlap (TMO) Earth Movers With Play-out Trace (EMT) Generalization by Trace Uniqueness (TGDU) Structural Simplicity by Edge Count (SSEDC) Structural Simplicity by Entity Count (SSENC) Structural Simplicity incl. Stochastic (SSS) Trace Ratio Gower 4 (TRG4) Trace Ratio Gower 3 (TRG3) Trace Ratio Gower 2 (TRG2) Activity Ratio Gower (ARG) Earth Movers Truncated (tEMSC0.8) Entropy Precision (H P) Entropy Recall (H F) Trace Overlap Ratio (TOR) Generalization by Trace Floor 1 (TGF1) Generalization by Trace Floor 10 (TGF10)





# Three Empirical Dimensions

### Adhesion

How little effort is required to transform one stochastic language into another

#### Entropy

The amount of information in a system

In this case, the combination of log and model

### Simplicity Structural simplicity of the model

Play-out Entropy Project Fitness (HJFT) Generalization by Trace Floor 5 (TGE5) Play-out Entropy Project Precision (HJPT) Play-out Entropy Intersection Fitness (HIFT) Play-out Entropy Intersection Precision (HIPT) Trace Probability Mass Overlap (TMO) Earth Movers With Play-out Trace (EMT) Generalization by Trace Uniqueness (TGDU) Structural Simplicity by Edge Count (SSEDC) Structural Simplicity by Entity Count (SSENC) Structural Simplicity incl. Stochastic (SSS) Trace Ratio Gower 4 (TRG4) Trace Ratio Gower 3 (TRG3) Trace Ratio Gower 2 (TRG2) Activity Ratio Gower (ARG) Farth Movers Truncated (tEMSC0.8) Entropy Precision (H P) Entropy Recall (H E) Trace Overlap Ratio (TOR) Generalization by Trace Floor 1 (TGF1) Generalization by Trace Floor 10 (TGF10) Scaled contributions

per component



Adhesion entropy simplicity



Adhesion + entropy simplicity



Adhesion + entropy + simplicity



Adhesion + entropy + simplicity +



Adhesion entropy simplicity



Adhesion + entropy simplicity



Adhesion + entropy - simplicity



Adhesion + entropy - simplicity +



Adhesion entropy simplicity



Adhesion  $\sim$  entropy simplicity



Adhesion  $\sim$  entropy  $\sim$  simplicity



Adhesion  $\sim$  entropy  $\sim$  simplicity -



Adhesion entropy simplicity



Adhesion - entropy simplicity



Adhesion - entropy - simplicity



Adhesion - entropy - simplicity  $\sim$ 



 $\blacktriangleright$  First models are process trees  $\rightarrow$  representational bias

- $\blacktriangleright$  First models are process trees  $\rightarrow$  representational bias
- SETM evolutionary fitness function may tend to correlate measures

- $\blacktriangleright$  First models are process trees  $\rightarrow$  representational bias
- SETM evolutionary fitness function may tend to correlate measures
  - Robustness tests excluding SETM still show the effect, though

- $\blacktriangleright$  First models are process trees  $\rightarrow$  representational bias
- SETM evolutionary fitness function may tend to correlate measures
  - Robustness tests excluding SETM still show the effect, though
- Largest log 200 000 traces

- Three empirically derived dimensions
- Focus on empirical and orthogonality
- Other measures and principles may be non-orthogonal but still useful, eg recall and precision entropy measures
- Future work
  - Theoretical grounded measures for these dimensions
  - Further tests