# Project Proposal: Prediction by Compression

Lasse Blaauwbroek

Czech Institute for Informatics, Robotics and Cybernetics
Czech Technical University in Prague

AITP 2018

March 30, 2018

Compressor $C$ such that $C(s)$ is the length of the compression of $s$

[Cilibrasi and Vitanyi 2003], [Li et al. 2004]

Compressor $C$ such that $C(s)$ is the length of the compression of $s$

$s$ and $t$ share all information $\implies C(st) \approx C(s) + b$

$s$ and $t$ share no information $\implies C(st) \approx C(s) + C(t)$

[Cilibrasi and Vitanyi 2003], [Li et al. 2004]

Compressor $C$ such that $C(s)$ is the length of the compression of $s$

$$s \text{ and } t \text{ share all information} \implies C(st) \approx C(s) + b$$
$$s \text{ and } t \text{ share no information} \implies C(st) \approx C(s) + C(t)$$

$$NCD_C(s, t) = \frac{C(st) - \min(C(s), C(t))}{\max(C(s), C(t))}$$

[Cilibrasi and Vitanyi 2003], [Li et al. 2004]

Compressor $C$ such that $C(s)$ is the length of the compression of $s$

$s$ and $t$ share all information $\implies C(st) \approx C(s) + b$

$s$ and $t$ share no information $\implies C(st) \approx C(s) + C(t)$

$$NCD_C(s, t) = \frac{C(st) - \min(C(s), C(t))}{\max(C(s), C(t))}$$

Under reasonable conditions for $C$, $NCD_c$ approximates a metric

[Cilibrasi and Vitanyi 2003], [Li et al. 2004]

Let $P$ be the set of valid programs for programming language $L$

[Cilibrasi and Vitanyi 2003], [Li et al. 2004]

Let $P$ be the set of valid programs for programming language $L$
Kolmogorov complexity $K$:

$$K(s) = \underset{p \in P \land L(p) = s}{\arg\min} |p|$$

[Cilibrasi and Vitanyi 2003], [Li et al. 2004]

Let $P$ be the set of valid programs for programming language $L$
Kolmogorov complexity $K$:

$$K(s) = \underset{p \in P \wedge L(p) = s}{\arg\min} |p|$$

$$NCD_K(s, t) = \frac{K(st) - \min(K(s), K(t))}{\max(K(s), K(t))}$$

$NCD_K$ is the distance metric:

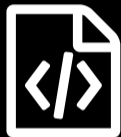$$\forall_{d,s,t} \; \text{computable}(d) \Rightarrow NCD_K(s, t) \leq d(s, t)$$

[Cilibrasi and Vitanyi 2003], [Li et al. 2004]

No domain-specific knowledge necessary!

No domain-specific knowledge necessary!

No domain-specific knowledge necessary!
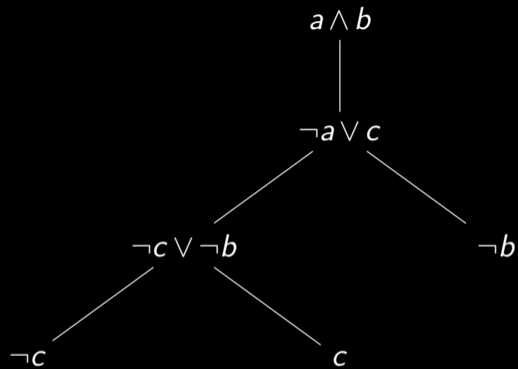
Problem: Mathematical statements are short

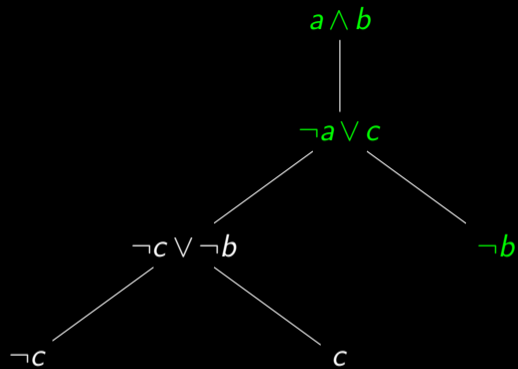Problem: Mathematical statements are short

# Compression: Prediction by Partial Matching
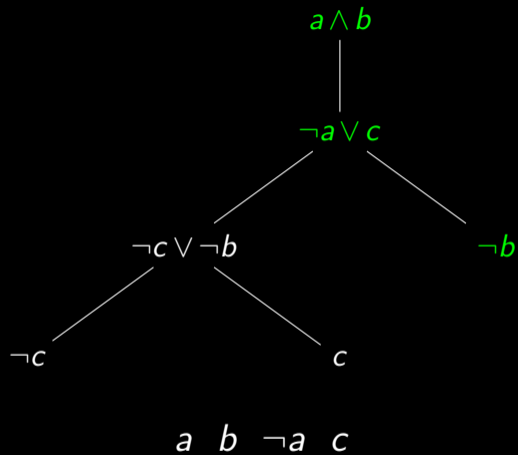
Problem: Mathematical statements are short

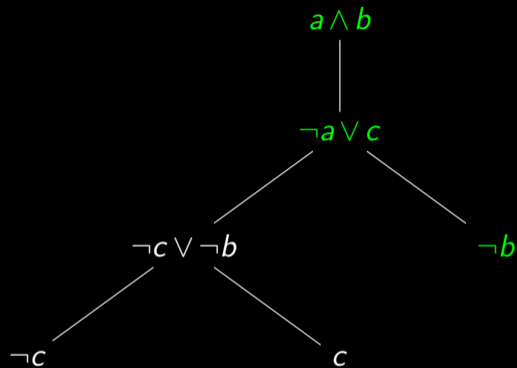# Compression: Prediction by Partial Matching

## Compress entire proof states

$a \wedge b$

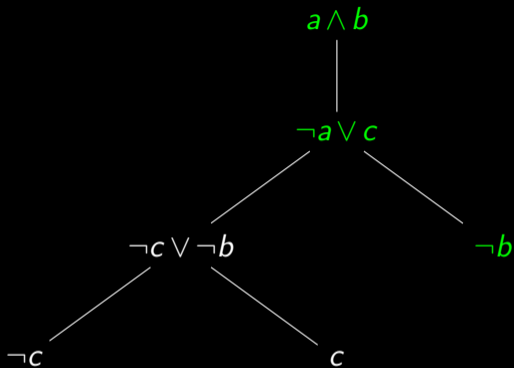$\neg a \vee c$

$\neg c \vee \neg b$                    $\neg b$

$\neg c$                    $c$

[CKaliszyk, Urban and Vyskoci 2015]

$a \wedge b$

$\neg a \vee c$

$\neg c \vee \neg b$

$\neg b$

$\neg c$

$c$

[CKaliszyk, Urban and Vyskoci 2015]

$a \wedge b$

$\neg a \vee c$

$\neg c \vee \neg b$　　　　$\neg b$

$\neg c$　　　　$c$

$a$　$b$　$\neg a$　$c$

[CKaliszyk, Urban and Vyskoci 2015]

$a \wedge b$

$\neg a \vee c$

$\neg c \vee \neg b$            $\neg b$

$\neg c$                $c$

$a \quad b \quad \neg a \quad c$

"$a \wedge b \neg a \vee c \neg bab \neg ac$"

[CKaliszyk, Urban and Vyskoci 2015]

$$a \wedge b$$

$$\neg a \vee c$$

$$\neg c \vee \neg b \qquad \neg b$$

$$\neg c \qquad c$$

$$a \quad b \quad \neg a \quad c$$

"$a \wedge b \neg a \vee c \neg bab \neg ac$"

Database

$\Longleftrightarrow$

[CKaliszyk, Urban and Vyskoci 2015]

About 30-40 compressions per second
No vector space: $n$ compressions per prediction

About 30-40 compressions per second
No vector space: $n$ compressions per prediction

Idea: Impose structure through an $n$-dimensional lattice

$$S_n = \{X \subseteq S \mid |X| = n\}$$

$$\text{out}(s) = \underset{X \in S_n}{\arg\max} \frac{\sum\limits_{t,u \in X} NCD(t, u)}{\sum\limits_{t \in X} NCD(s, t)}$$

## Pros

▷ No domain-specific knowledge required
▷ Predictions are competitive
▷ Robust against different representations of proof states

<center>Pros</center>

▷ No domain-specific knowledge required
▷ Predictions are competitive
▷ Robust against different representations of proof states

<center>Cons</center>

▷ Relatively slow
▷ No vector space

<div align="center">Pros</div>

▷ No domain-specific knowledge required
▷ Predictions are competitive
▷ Robust against different representations of proof states

<div align="center">Cons</div>

▷ Relatively slow
▷ No vector space

<div align="center">Ideas</div>

▷ Adapt the PPM compressor for tree-structures
▷ Impose a $n$-dimensional lattice on the data

<center>Pros</center>

▷ No domain-specific knowledge required
▷ Predictions are competitive
▷ Robust against different representations of proof states

<center>Cons</center>

▷ Relatively slow
▷ No vector space

<center>Ideas</center>

▷ Adapt the PPM compressor for tree-structures
▷ Impose a $n$-dimensional lattice on the data

<center>?</center>