# Learning in Pessiland
# via
# Inductive Inference

Shuichi Hirahara

NII

Mikito Nanashima

TokyoTech

**Online Complexity Seminar**

# Backgrounds

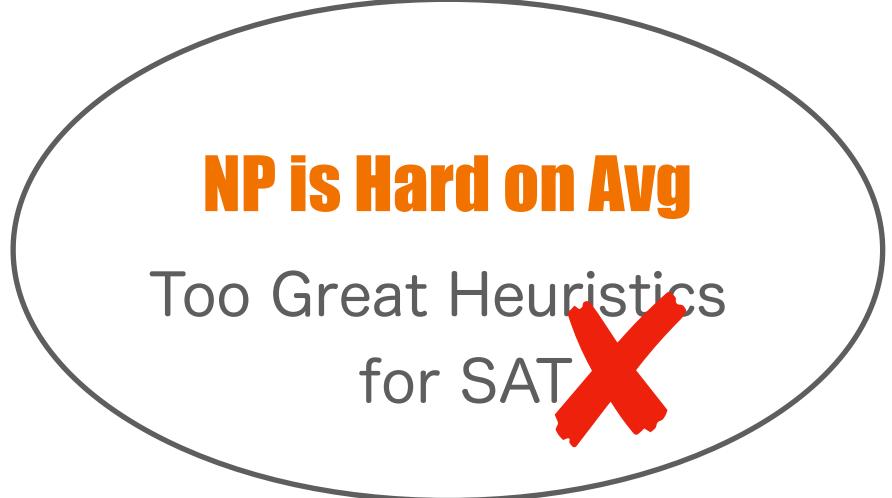# Our Results

# Proof Techniques

# Conclusion

# Backgrounds

Our Results

Proof Techniques

Conclusion

# Pessiland

NP is Hard on Avg

Too Great Heuristics for SAT ❌

No One-Way Functions

(Classic) Crypto ❌

**Suppose Our World is Unfortunately Pessiland.**
**What Can We Do?**

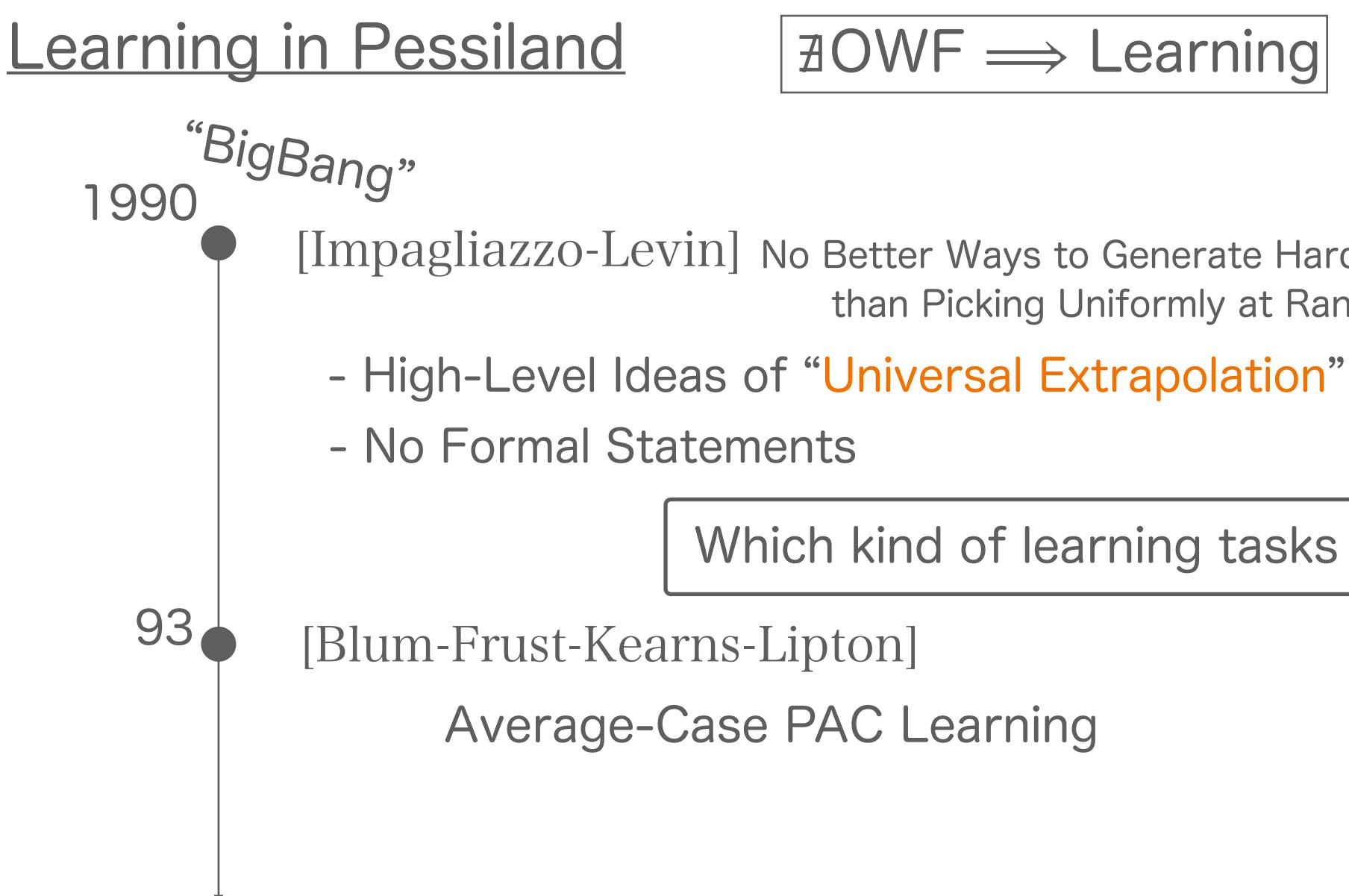Average-Case Inverter for Poly-Time Functions
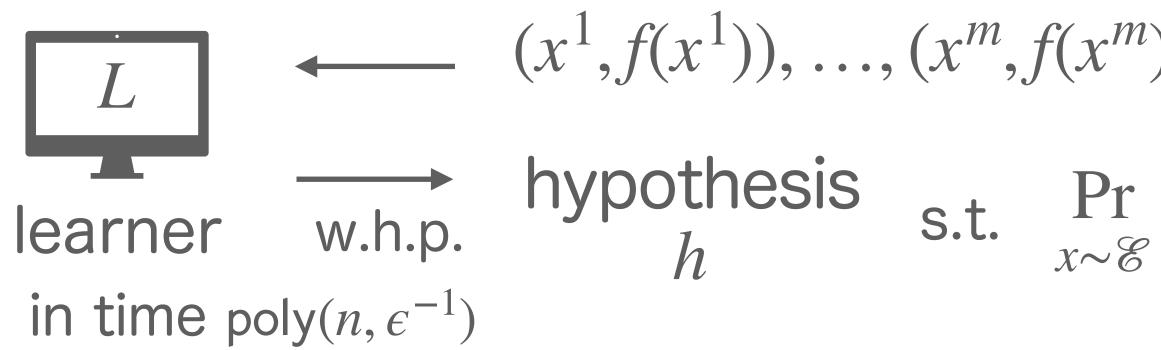
**Algorithms in Pessiland**

Hardness

# Learning in Pessiland

$\boxed{\nexists \text{OWF} \implies \text{Learning}}$

*"BigBang"*

**1990**

[Impagliazzo-Levin] No Better Ways to Generate Hard NP Instances
than Picking Uniformly at Random

- High-Level Ideas of "Universal Extrapolation"

- No Formal Statements

Which kind of learning tasks can be done ?

**93**

[Blum-Frust-Kearns-Lipton]

Average-Case PAC Learning

# PAC Learning in Pessiland

$\nexists$OWF$\implies$ Avg PAC learning

a concept class $\mathscr{C} = \{\mathscr{C}_n\}_{n\in\mathbb{N}}$    $\mathscr{C}_n \subseteq \{f\colon \{0,1\}^n \to \{0,1\}\}$

an example distribution   $\mathscr{E} = \{\mathscr{E}_n\}_{n\in\mathbb{N}}$    $\mathscr{E}_n$ is over $\{0,1\}^n$

a distribution over functions   $\mathscr{F} = \{\mathscr{F}_n\}_{n\in\mathbb{N}}$   $\mathscr{F}_n$ is over $\mathscr{C}_n$

$$\boxed{L} \longleftarrow (x^1, f(x^1)), \ldots, (x^m, f(x^m)) \qquad \begin{array}{l} f \sim \mathscr{F}_n \\ x^1, \ldots, x^m \sim \mathscr{E}_n \end{array}$$

learner   $\xrightarrow{\text{w.h.p.}}$   hypothesis $h$   s.t.   $\Pr_{x\sim\mathscr{E}}[h(x) \neq f(x)] \leq \epsilon$

in time $\mathrm{poly}(n, \epsilon^{-1})$

[BFKL93]   $\mathscr{C}$ : efficiently evaluatable,   $\mathscr{E}, \mathscr{F}$ : samplable

# Learning in Pessiland

$\boxed{\nexists \text{OWF} \implies \text{Learning}}$

*"BigBang"*

**1990**

[Impagliazzo-Levin] No Better Ways to Generate Hard NP Instances
than Picking Uniformly at Random
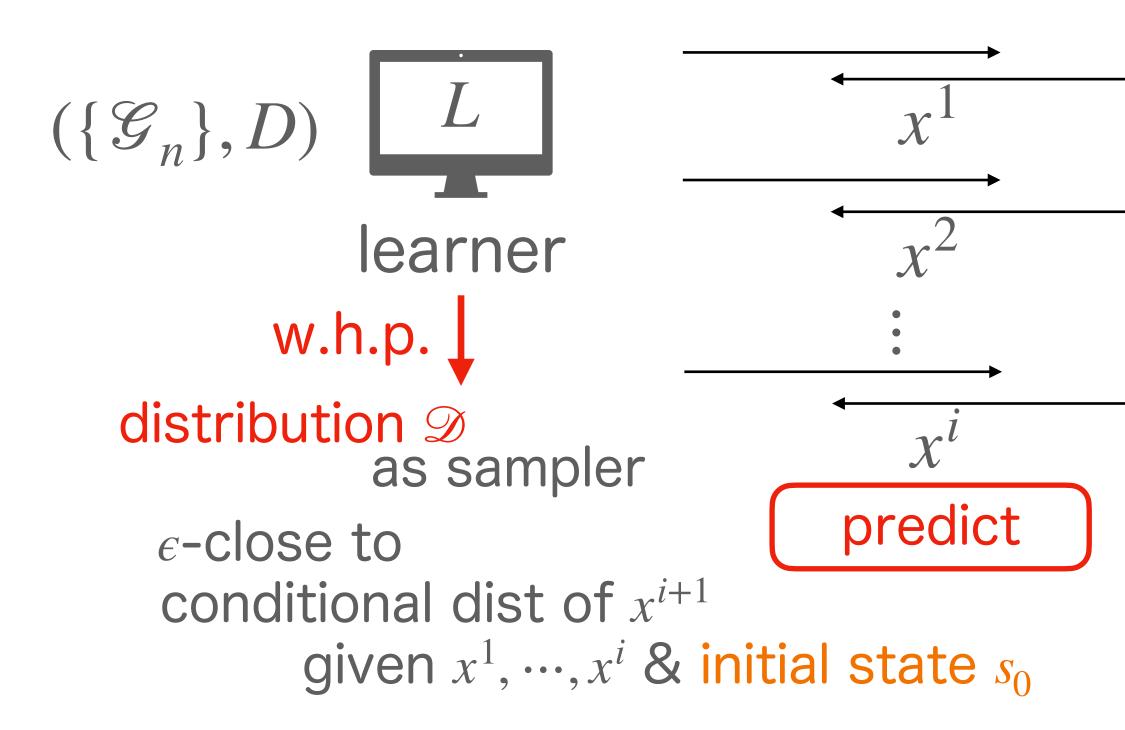
- High-Level Ideas of "Universal Extrapolation"

- No Formal Statements

**93** [Blum-Frust-Kearns-Lipton]

Average-Case PAC Learning

**2006** [Naor-Rothblum]

Learning Adaptively Changing Distributions (ACDs)

# Distributional Learning in Pessiland

[NR06]   $\not\exists$OWF$\implies$ Learning known ACDs

ACD  $(\{\mathscr{G}_n\}, D)$

$\mathscr{G}$: samplable   $D$ : poly-time sampler

$(\{\mathscr{G}_n\}, D)$

$L$

learner

w.h.p. ↓

distribution $\mathscr{D}$

as sampler

$x^1$

$x^2$

⋮

$x^i$

predict

$\epsilon$-close to
conditional dist of $x^{i+1}$
given $x^1, \cdots, x^i$ & initial state $s_0$

internal state $s$  $= s_0$

Initialization:  $s_0 \sim \mathscr{G}_n$

initial state  **target**

seed
$r \sim \{0,1\}^{\mathsf{poly}(n)}$

sample
$x \in \{0,1\}*$

$D$

state   $s$

$s'$

$s := s'$  next state

[NR06]   $L$ uses the knowledge of $\{\mathscr{G}_n\}$ and $D$

# Why is the knowledge of distributions important?

Critical Cases

$(\{\mathscr{G}_n\}, D)$ ✗

ACD $(\{\mathscr{G}_n\}, D)$

$\mathscr{G}$: samplable    $D$ : poly-time sampler



$L$

learner

w.h.p.

distribution $\mathscr{D}$
as sampler

$\epsilon$-close to
conditional dist of $x^{i+1}$
given $x^1, \cdots, x^i$ & initial state $s_0$

internal state $s = s_0$

Initialization:  $s_0 \sim \mathscr{G}_n$

initial state    **target**

specific
$i \in \mathbb{N}$

seed
$r \sim \{0,1\}^{\mathsf{poly}(n)}$

$x^1$

$x^2$

$\vdots$

$x^i$

$s_0$ ✗

use $s_0$

state $s$

$x \in \{0,1\}*$

$D$

$s'$

$s := s'$    next state

[NR06]   $L$ uses the knowledge of $\{\mathscr{G}_n\}$ and $D$

to test whether the next sample is predictable

# Improved Learning in Pessiland ?

[BFKL]   examples and target function are separately selected

➡   Q.  Average-Case PAC Learning on Joint Distribution?

$$\mathscr{D} \sim \mathscr{G} \qquad (x, b) \sim \mathscr{D}$$
↑
description of
joint distribution

Separated (BFKL)      Joint
proper learning  → DistNP-hard
                      [Pitt-Valiant]

Q.  Agnostic Learning ?    $\Pr_{(x,b)} [h(x) \neq b] \leq \min_{f \in \mathscr{C}} \Pr_{(x,b)} [f(x) \neq b] + \epsilon$

[NR]   learning known ACDs

➡   Q.  Learning unknown ACDs ?

Backgrounds

Our Results

Proof Techniques

Conclusion

# Our Contribution

Resolve   Q. Joint Distribution?   Q. Learning unknown ACDs ?

Q. Agnostic Learning ?

1990

● Impagliazzo-Levin, No Better Ways to Generate Hard NP Instances
than Picking Uniformly at Random

- High-Level Ideas of "Universal Extrapolation"   revisit

- No Formal Statements

93 ● Blum-Frust-Kearns-Lipton

Average-Case PAC Learning

Unified &
Simplified &
Improved

2006

● Naor-Rothblum

Learning Adaptively Changing Distributions (ACDs)

# In this talk...

1. $\not\exists$ Infinitely Often OWF

"Adversaries can invert functions for all sufficiently large parameters"

$\not\exists$ (standard) OWF $\longrightarrow$ infinitely many size $n$

accuracy, confidence $\leq 1/\mathrm{poly}(n)$

fixed as poly-time functions in $n$

2. No details for the choices of parameters

I do not discuss confidence parameters so much

minimize $\mathrm{K}^{\text{your\_time}}$(the paper | this talk)

# Learning ACDs in Pessiland

Q. Learning unknown ACDs ?

---

**Thm. 1**

$\nexists$OWF iff $\exists$ poly-time learner for all (unknown) ACDs $(\mathscr{G}, D)$

with sample complexity $\underline{O(s\epsilon^{-2})}$ for $s$-bit initial states and accuracy $\epsilon$

---

improved from $O(s\epsilon^{-4})$ [NR06]

How can we avoid the critical cases?



$L$

learner

for unknown $(G, D)$

$x^1$

$x^2$

$\vdots$

$x^i$

$r \searrow$
$\phantom{r}\; D \phantom{=}$
$s \nearrow$

$\searrow x$
$\nearrow s'$

succeeds
at almost all steps

choose a prediction stage
uniformly at random

# Agnostic Learning in Pessiland

Q. Joint Distribution?      Q. Agnostic Learning ?

---

**Thm. 2**

$\nexists$OWF iff $\exists$ poly-time agnostic learner for $\mathscr{F} = \{f\colon \{0,1\}^n \to \{0,1\}^{\text{poly}(n)}\}$

(with 0-1 loss)

(= learning by information theoretically optimal hypothesis)

on avg under a joint dist $\mathscr{D}$ on samples, where $\mathscr{D} \sim \mathscr{G}$, $\mathscr{G}$ is samplable

with sample complexity $O(s\epsilon^{-2})$    (for $s = |\mathscr{D}|$, accuracy $\epsilon$)

---

optimal in general

[Previous]

PAC

Separated Distributions

Binary Labels

▶

[Ours]

Agnostic

Joint Distributions

Multi Labels

Improper learning
(General Hypothesis)

# Improved Learning in Pessiland

$\nexists$OWF $\iff$ Worst-Case Learning

<span style="color:red">exp-time in Computational Depth of Secrets</span>

$U$: universal TM

$$\mathrm{K}(x) = \min\{p \in \mathbb{N} : \exists \Pi \in \{0,1\}^p \text{ s.t. } U(\Pi) = x\}$$

$\mathrm{Q}^t := $ dist. of $U(w)$ executed in $t$ steps for $w \sim \{0,1\}^t$

$\mathrm{q}^t(x) := -\log \Pr[x \sim \mathrm{Q}^t]$         $\mathrm{q}^t(x) \approx \mathrm{pK}^{t'}(x)$  introduced in [GKLO22]

$\mathrm{cd}^t(x) := \mathrm{q}^t(x) - \mathrm{K}(x)$

<span style="color:orange">$\approx \mathrm{pK}^{t'}(x) - \mathrm{pK}^{\infty}(x)$</span>

$$\boxed{\begin{array}{ll} \mathrm{pK}^{\mathrm{poly}(t)}(x) \lesssim \mathrm{q}^t(x) & \text{Optimal coding } {}_{\text{[LOZ22]}} \\ \mathrm{pK}^t(x) \gtrsim \mathrm{q}^{\mathrm{poly}(t)}(x) & \text{Domination} \end{array}}$$

$$\boxed{\begin{array}{c} \text{for any samplable distribution } \mathcal{D} = \{\mathcal{D}_n\} \\ \mathrm{cd}^{\mathrm{poly}}(x) = O(\log n) \text{ w.h.p. } x \sim \mathcal{D}_n \end{array}}$$

$$\boxed{\begin{array}{c} \textbf{Slow Growth Law} \\ x \longrightarrow \boxed{M:\mathsf{PPT}} \longrightarrow x' \\ \mathrm{cd}^{O(t+time_M)}(x') \lesssim \mathrm{cd}^t(x) \text{ w.h.p.} \end{array}}$$

# Improved Learning in Pessiland

$\nexists$OWF $\iff$ Worst-Case Learning

exp-time in Computational Depth of Secrets

---

**Thm. 3**

The following are equivalent:

1. $\nexists$ OWF

2. learning unknown ACDs in time $\mathrm{poly}(n, t, \epsilon^{-1}, 2^{\mathrm{cd}^t(s_0)})$, where $s_0$ is initial state

   (worst case on initial states)

3. agnostic learning on unknown joint dists $\mathscr{D}$ over samples

   in time $\mathrm{poly}(n, t, \epsilon^{-1}, 2^{\mathrm{cd}^t(|\mathscr{D}|)})$

   (worst case on joint distributions over samples)

---

**Note** Thm1 & 2 are implied by Thm3

Backgrounds

Our Results

Proof Techniques

Conclusion

# Our Approach

**Step I**  State "Universal Extrapolation" <span style="color:red">formally</span>

**Step II**  Translate "Universal Extrapolation" into Learning

# Our Approach

**Step I**     State "Universal Extrapolation" <span style="color:red">formally</span>

**Step II**    Translate "Universal Extrapolation" into Learning

# Formulating Universal Extrapolation

**Our Proposal**   Universal Extrapolation = Extrapolation under $Q^t$

(Time-Bounded Universal Distribution)

**Notation**   distribution $\mathscr{D}$   prefix $x \in \{0,1\}^*$   $k \in \mathbb{N}$

$\text{Next}_k(x; \mathscr{D})$   = distribution of k bits following $x$ w.r.t. $\mathscr{D}$

prefix $x \in \{0,1\}^*$



paramaters
$k \in \mathbb{N}\ t \in \mathbb{N}\ \epsilon \in (0,1)$

$y \in \{0,1\}^{\leq k}$  $\approx$  $\text{Next}_k(x; Q^t)$

within statistical distance $\epsilon$

**Lemma** (informal)   $\not\exists$ (io)OWF $\implies$  $\exists$UE  that works **in worst case on $x$**

in time $\text{poly}(|x|, k, t, \epsilon^{-1}, 2^{\text{cd}^t(x)})$

(UE is given $2^\alpha$ (in unary) and works for every $x$ with $\text{cd}^t(x) < \alpha$)

**Lemma** (informal)    $\nexists$ (io)OWF $\implies$ $\exists$UE that works **in worst case on** $x$

in time $\mathrm{poly}(|x|, k, t, \epsilon^{-1}, 2^{\mathrm{cd}^t(x)})$

## Distributional Inverting

Inverting

Distributional Inverting

some $x'$

$$f(x) \longrightarrow \boxed{Inv} \longrightarrow f(x') = f(x)$$

$(x \sim \{0,1\}^*)$

some $x'$

$$f(x) \longrightarrow \boxed{DInv} \longrightarrow f(x') = f(x)$$

$(x \sim \{0,1\}^*)$

simulate unif sampling
from $\{x' : f(x') = f(x)\}$

**Thm** [IL89]

$\nexists$ (io)OWF $\implies$ for every poly-time function $f = \{f_n\}$

$\exists DInv$ : PPT s.t. $\forall n, \epsilon^{-1}, \delta^{-1} \in \mathbb{N}$

$$\Delta_{TV}\left(DInv(f_n(x); 1^n, 1^{\epsilon^{-1}}, 1^{\delta^{-1}}), \text{Unif over } f_n^{-1}(f_n(x))\right) \leq \epsilon$$

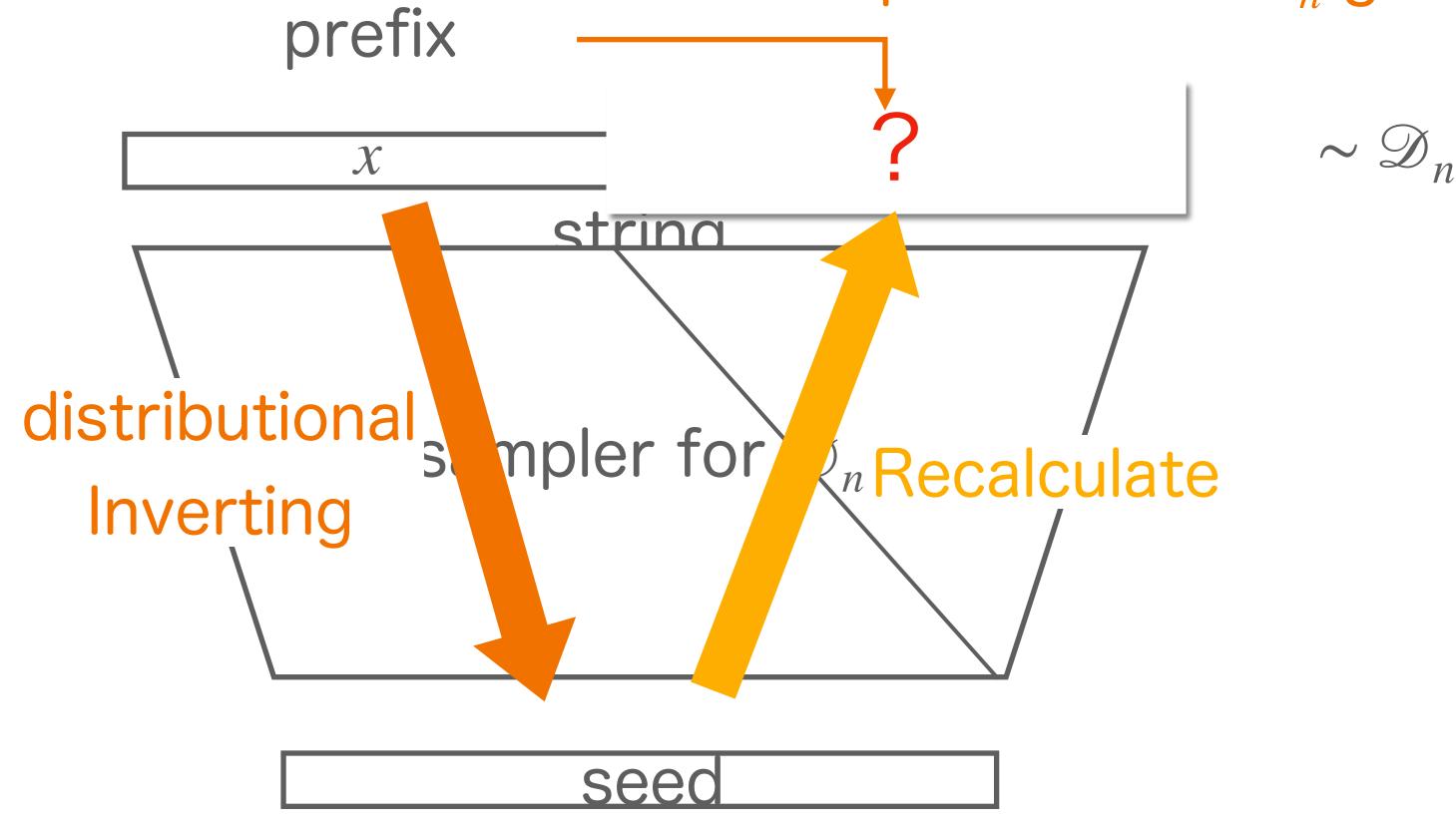w.p. $\geq 1 - \delta$ over $x \sim \{0,1\}^{\mathrm{poly}(n)}$

**Lemma** (informal)     $\nexists$ (io)OWF $\implies$     $\exists$UE  that works **in worst case on** $x$

in time $\mathrm{poly}(|x|, k, t, \epsilon^{-1}, 2^{\mathrm{cd}^t(x)})$

**Distributional Inverting $\rightarrow$ Distribution-Specific Extrapolation**   [Ost91, OW93, NR06, . . .]

$\mathscr{D} = \{\mathscr{D}_n\}$ samplable distribution

extrapolate w.r.t $\mathscr{D}_n$ given $x$

prefix

$x$

$\sim \mathscr{D}_n$

?

string

distributional
Inverting

sampler for $\mathscr{D}_n$ Recalculate

seed

**Lemma** (informal) $\quad \nexists$ (io)OWF $\implies$ $\exists$UE that works **in worst case on** $x$ in time $\text{poly}(|x|, k, t, \epsilon^{-1}, \textcolor{red}{2^{\text{cd}^t(x)}})$

## Distribution-Specific Extrapolation for $Q^t$

$\exists \text{UE}' : \text{PPT s.t. } \forall t, \ell, k, \epsilon^{-1}, \delta^{-1} \in \mathbb{N}$

$$\Delta_{TV}\left(\text{UE}'(x), \text{Next}_k(x; Q^t)\right) \leq \epsilon \quad \text{w.p. } \geq 1 - \delta \text{ over } x \sim Q^t_{\leq \ell}$$

the first $\ell$ bits of $Q^t$

Based on [AF09] $\quad E := E_{t,\ell,k,\epsilon,\delta} = \left\{ x : \Delta_{TV}\left(\text{UE}'(x), \text{Next}_k(x; Q^t)\right) > \epsilon \right\}$ error set

$\delta := 2^{-\alpha}$ Goal: $x \in E \implies \text{cd}^t(x) \geq \alpha - O(\log t\ell k\epsilon^{-1}\alpha)$

$$2^{-\alpha} \geq \sum_{x \in E} \Pr[x \sim Q^t_{\leq \ell}] \geq \sum_{x \in E} \Pr[x \sim Q^t] \qquad \sum_{x \in E} \underline{2^\alpha \Pr[x \sim Q^t]} \leq 1$$

computable given $\text{UE}'$ & parameters

inefficiently computable distribution $\{\mathscr{E}_{t,\ell,k,\epsilon,\delta}\}$ $\quad \forall x \in E \ \Pr[x \sim \mathscr{E}] = 2^\alpha \Pr[x \sim Q^t]$

$\forall x \in E \quad \text{K}(x) \leq -\log \Pr[x \sim \mathscr{E}] + O(\log t\ell k\epsilon\alpha)$ optimal coding

$$= -\alpha + \underbrace{(-\log \Pr[x \sim Q^t])}_{q^t(x)} + O(\log t\ell k\epsilon\alpha)$$

19/28

# Step 1: Summary

**Our Proposal**  Universal Extrapolation = Extrapolation under $Q^t$

(Time-Bounded Universal Distribution)

prefix $x \in \{0,1\}*$

paramaters
$k \in \mathbb{N}$  $t \in \mathbb{N}$  $\epsilon \in (0,1)$



$y \in \{0,1\}^k$  $\approx$  $\mathrm{Next}_k(x; Q^t)$

with statistical distance within $\epsilon$

**Lemma** (informal)  $\nexists$ (io)OWF  $\implies$  $\exists$UE  that works **in worst case on** $x$

in time $\mathrm{poly}(|x|, k, t, \epsilon^{-1}, 2^{\mathrm{cd}^t(x)})$

(UE is given $2^\alpha$ (in unary) and works for every $x$ with $\mathrm{cd}^t(x) < \alpha$)

**Proof**  Use Distributional Inverter for $Q^t$

**Q. How can we obtain learners (e.g., agnostic learners)?**

# Step2 UE → Learning

offline stream

efficiently &
randomly

$$x^1 \ y^1 \ x^2 \ y^2 \ \ldots \ x^m \ y^m \longleftarrow$$

secret info.
(sampler)
$\Pi \in \{0,1\}*$

possibly correlated

$$x^1 \ y^1 \ x^2 \ y^2 \ \ldots \ x^{i-1} \ y^{i-1} \ x^i \qquad \boxed{y^i \ \text{?}}$$

$x^1, y^1, x^2, y^2 \ldots \in \{0,1\}*$

$L$

**predict**

$x$ : advice information

$y$ : target

**EX.** learning ACDs

$\Pi \approx$ initial state   $x$ : empty   $y$ : sample
simulate the distribution of $y^i$

agnostic learning
(with 0-1 loss)

$\Pi \approx$ joint distribution over samples
$x$ : example   $y$ : label
answer the best possible $y^i$

# Cheating Learner

offline stream    efficiently & randomly

$x^1 \; y^1 \; x^2 \; y^2 \; \ldots \; x^m \; y^m$ ←——

secret info.
(sampler)
$\Pi \in \{0,1\}*$

possibly correlated

$x^1 \; y^1 \; x^2 \; y^2 \; \ldots \; x^{i-1} \; y^{i-1} \; x^i$    $\boxed{y^i \; ?}$

$x^1, y^1, x^2, y^2 \ldots \in \{0,1\}*$

$x$ : advice information

$y$ : target

$L^{\mathcal{O}}$

condition ↓

$\mathcal{O}$

$h^{\mathcal{O}}$

$y^i \; y^i \; y^i$ ←—— $\mathcal{O}$

w.r.t. $\Pi$

**EX.** learning ACDs

$\Pi \approx$ initial state    $x$ : empty    $y$ : sample

simulate the distribution of $y^i$

✅ by outputting $y^i \sim \mathcal{O}$

agnostic learning
(with 0-1 loss)

$\Pi \approx$ joint distribution over samples

$x$ : example    $y$ : label

answer the best possible $y^i$

✅ by collecting $y^i, y^i, \ldots \sim \mathcal{O}$

# UE → Learning

offline stream

efficiently & randomly

$x^1 \quad y^1 \quad x^2 \quad y^2 \quad ... \quad x^m \quad y^m$ ← secret info. (sampler) $\Pi \in \{0,1\}*$

possibly correlated

$x^1 \quad y^1 \quad x^2 \quad y^2 \quad ... \quad x^{i-1} \quad y^{i-1} \quad x^i$ $\boxed{y^i \; ?}$

$x^1, y^1, x^2, y^2 ... \in \{0,1\}*$

condition

$x$ : advice information

$L$ ✗

$y$ : target

$h$ ✗

$y^i \quad y^i \quad y^i$ ← ✗

w.r.t. $\Pi$

≋

w.h.p. when $i \sim [m]$

UE

for statistical distance within $\epsilon$

$m = O(|\Pi| \epsilon^{-2})$

**Q. Why ?**

learning ACDs ✅    agnostic learning (with 0-1 loss) ✅

sample complexity: linear in $|\Pi|$    time complexity: exp in $|\Pi|$

# Solomonoff's Inductive Inference [Sol64, LV19]

When extrapolating symbols, attach a higher probability for a more precise hypothesis, particularly, with an exponential rate on the description size

$q^{poly}(\cdot \mid \cdot)$ or $pK^{poly}(\cdot \mid \cdot)$

$(\because 2^{-q^{poly}} = Q^{poly})$  by conditional coding [HILNO23]

**Domination + Chain Rule for KL divergence**

efficiently & randomly

$$x^1 \quad y^1 \quad x^2 \quad y^2 \quad \dots \quad x^m \quad y^m \longleftarrow \boxed{\Pi \in \{0,1\}*}$$

When $t \gg \text{time}(\Pi)$

$$\Pr\left[x^1 y^1 \dots x^m y^m \sim Q^t\right] \geq 2^{-O(|\Pi|)} \Pr\left[x^1 y^1 \dots x^m y^m \sim \Pi\right] \quad (\textbf{Domination})$$

$$\log \frac{\Pr\left[x^1 y^1 \dots x^m y^m \sim \Pi\right]}{\Pr\left[x^1 y^1 \dots x^m y^m \sim Q^t\right]} \leq O(|\Pi|)$$

Taking the expectation over $x^1, y^1, \dots, x^m, y^m \sim \Pi$

$$\text{KL}\left(\Pi \,\middle\|\, Q^t\right) \leq O(|\Pi|)$$

# Solomonoff's Inductive Inference [Sol64, LV19]

$$\mathrm{KL}\left(\Pi \,\middle\|\, Q^t\right) \leq O(|\Pi|)$$

$$x^1 \; y^1 \; x^2 \; y^2 \; \cdots \; y^{i-1} \; x^i \; y^i \; \cdots \; x^m \; y^m$$

$X^1 \; Y^1 \qquad\qquad\qquad\qquad\qquad\qquad X^m \; Y^m$



**true distribution**

$\Pi \in \{0,1\}^*$

$\tilde{X}^1 \; \tilde{Y}^1 \qquad\qquad\qquad\qquad \tilde{X}^m \; \tilde{Y}^m$

$L \quad \square \; \square \quad x^1 \; y^1 \; x^2 \; y^2 \; \cdots \; y^{i-1} \; x^i \; y^i \; \cdots \; x^m \; y^m \quad \longleftarrow \quad Q^t$

UE — ignore the statistical error

$$O(|\Pi|) \geq \mathrm{KL}\left(X^1 Y^1 \cdots X^m Y^m \,\middle\|\, \tilde{X}^1 \tilde{Y}^1 \cdots \tilde{X}^m \tilde{Y}^m\right)$$

$$= \sum_{i=1}^{m} \mathrm{KL}\left((Y^i \,|\, X^1 Y^1 \cdots Y^{i-1} X^i) \,\middle\|\, (\tilde{Y}^i \,|\, \tilde{X}^1 \tilde{Y}^1 \cdots \tilde{Y}^{i-1} \tilde{X}^i)\right) \qquad \textbf{(Chain Rule)}$$

$$+ \sum_{i=1}^{m} \mathrm{KL}\left((X^i \,|\, X^1 Y^1 \cdots X^{i-1} Y^{i-1}) \,\middle\|\, (\tilde{X}^i \,|\, \tilde{X}^1 \tilde{Y}^1 \cdots \tilde{X}^{i-1} \tilde{Y}^{i-1})\right)$$

$$\geq \sum_{i=1}^{m} \mathrm{KL}\left((Y^i \,|\, X^1 Y^1 \cdots Y^{i-1} X^i) \,\middle\|\, (\tilde{Y}^i \,|\, \tilde{X}^1 \tilde{Y}^1 \cdots \tilde{Y}^{i-1} \tilde{X}^i)\right)$$

$\mathcal{O}_i \qquad\qquad \mathrm{UE}_i$

# Solomonoff's Inductive Inference [Sol64, LV19]

$$O(|\Pi|) \geq \sum_{i=1}^{m} \text{KL}(\ \boxed{i}\ \|\ \text{UE}_i\ )$$

$$\frac{O(|\Pi|)}{m} \geq \text{E}_{i\sim[m]}[\text{KL}(\ \boxed{i}\ \|\ \text{UE}_i\ )]$$

$$m \gg \frac{|\Pi|}{\epsilon^2} \implies \text{E}_{i\sim[m]}[\text{KL}(\ \boxed{i}\ \|\ \text{UE}_i\ )] \ll \epsilon^2$$

(**Markov**)　　$\text{KL}(\ \boxed{i}\ \|\ \text{UE}_i\ ) \ll \epsilon^2$　w.h.p. over $i \sim [m]$

(**Pinsker**)　　$\Delta_{TV}(\ \boxed{i}\ ,\ \text{UE}_i\ ) \leq \epsilon$　w.h.p. over $i \sim [m]$

# Remarks

$$x^1 \quad y^1 \quad x^2 \quad y^2 \quad \ldots \quad x^{i-1} \quad y^{i-1} \quad x^i \qquad \boxed{y^i \ \textbf{?}}$$

condition

$L$ ✗

$h$ ✗

$y^i \quad y^i \quad y^i \quad \longleftarrow \quad$ ✗ ✗ $\quad\approx\quad$ UE

when $i \sim [m]$   for stat dist within $\epsilon$

$$m = O(|\Pi|\epsilon^{-2})$$

Tasks poly-time cheating learners can do
= Tasks poly-time learners can do with UE

**Not sample optimal in agnostic learning**

**Why?**    accuracy $\longrightarrow$ query complexity of cheating learner

for optimal sample complexity

Extend **universal prediction** [MF98] to **computational cases**
**(statistical cases)**

Backgrounds


Our Results


Proof Techniques


Conclusion

# Summary

$$\nexists \text{OWF} \implies \text{Learning}$$



prefix
$x \in \{0,1\}*$ → UE → next $k$ bits under $Q^t$ → $L$ ⇄ UE

$x^1 \quad y^1 \quad \dots \quad x^{i-1} \quad y^{i-1} \quad x^i$ ?

---

**Thm.**

The following are equivalent:

1. $\nexists$ OWF

2. learning unknown ACDs in time $\text{poly}(t, \epsilon^{-1}, 2^{\text{cd}^t(s_0)})$, where $s_0$ is initial state

3. agnostic learning on unknown joint dists $\mathscr{D}$ over samples

   in time $\text{poly}(t, \epsilon^{-1}, 2^{\text{cd}^t(|\mathscr{D}|)})$

Q. weakest assumption for learning in time $\text{poly(cd)}$? AIOWF or HSG or $\cdots$?