

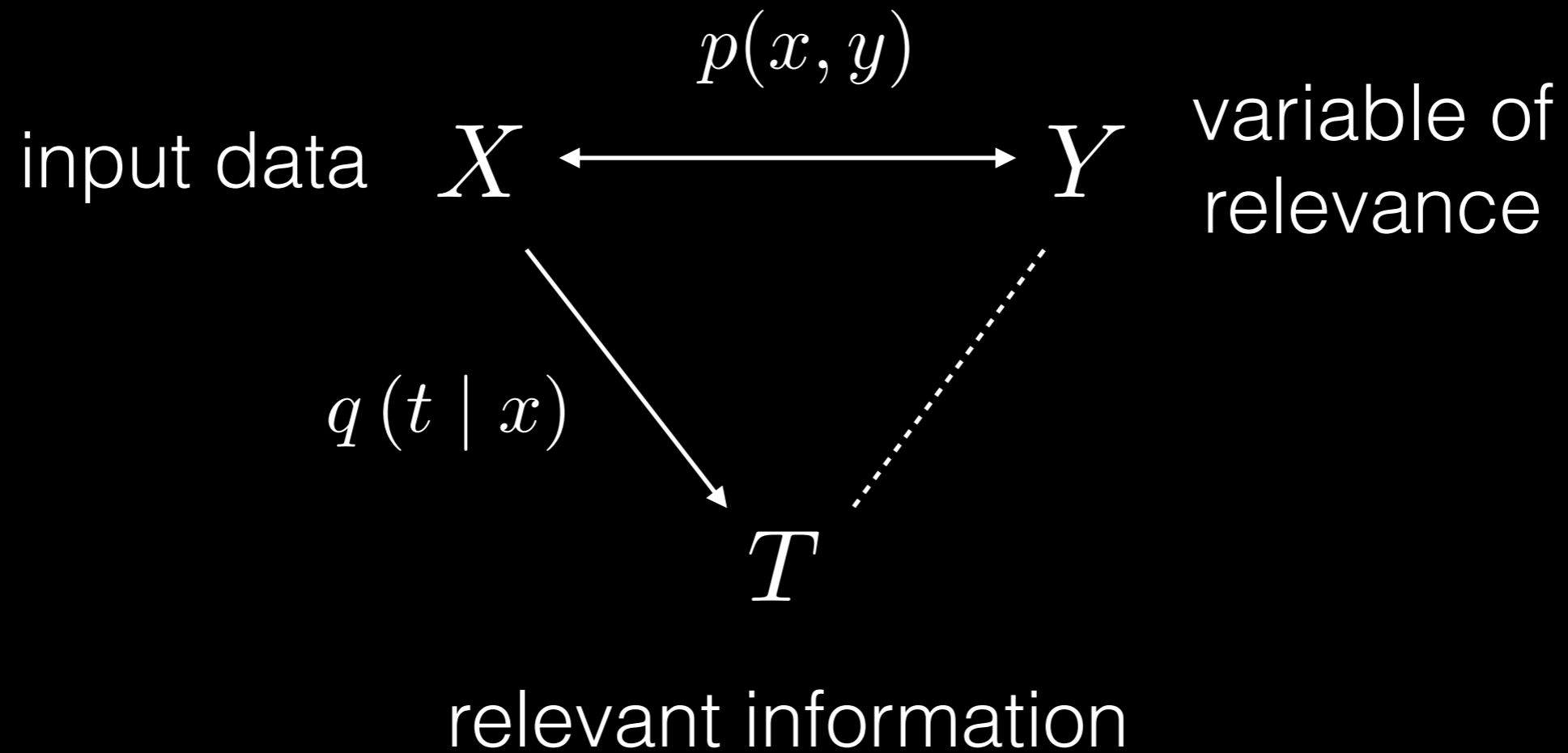
The deterministic information bottleneck

DJ Strouse
Princeton University

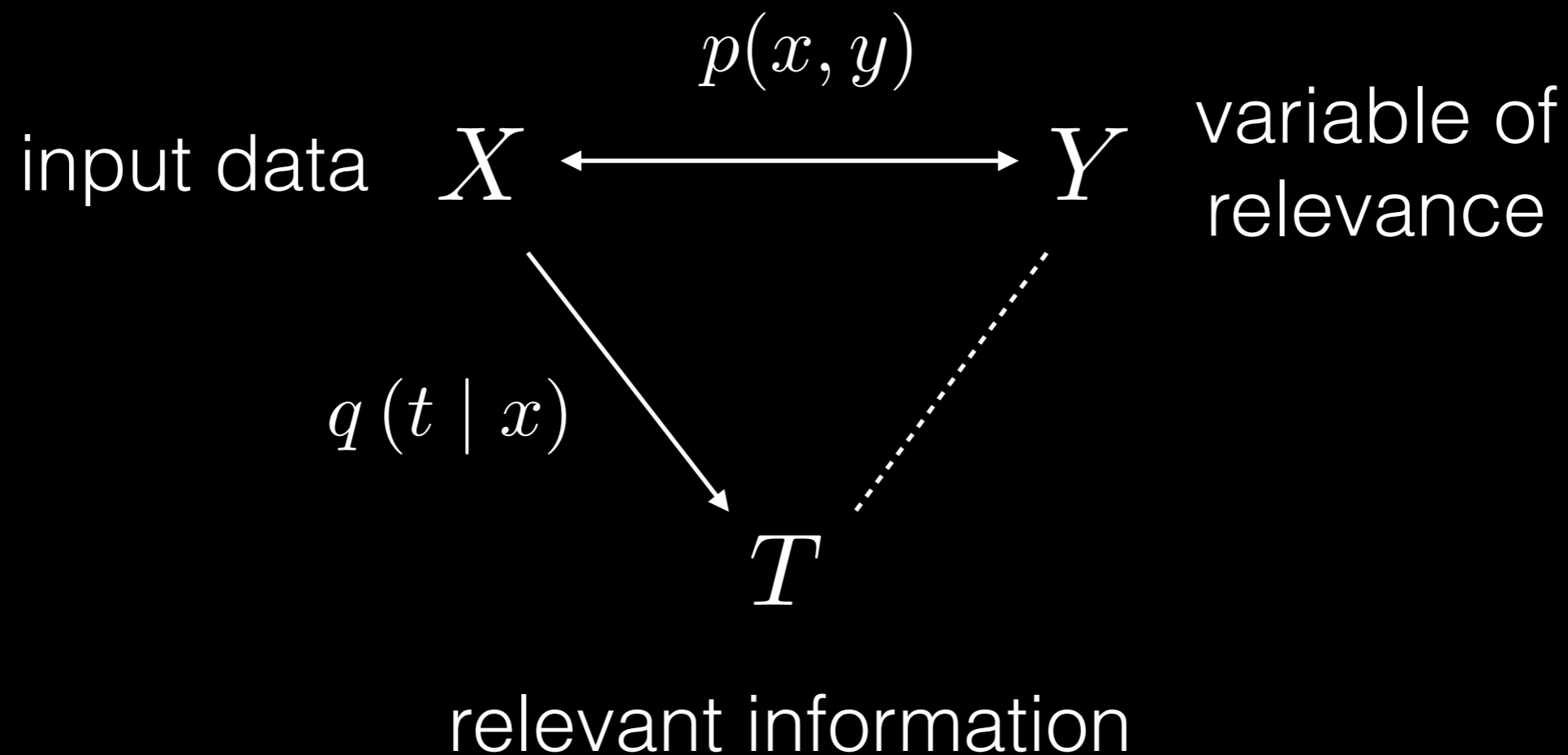
David Schwab
Northwestern University

DOE CSGF Program Review 2016
Washington, DC

Information bottleneck (IB)



Information bottleneck (IB)



statistics: soft sufficient statistic

info theory: lossy compression, distortion \sim relevance

machine learning: maximally informative clustering

IB examples

X

T

Y

user segmentation	demographics & past behavior	cluster ID	future purchase/ click behavior
------------------------------	---------------------------------	------------	------------------------------------

IB examples

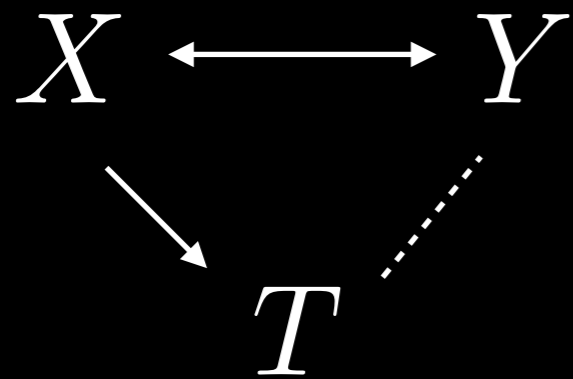
X

T

Y

user segmentation	demographics & past behavior	cluster ID	future purchase/ click behavior
human attention & memory	sensory input	neural activity/ synaptic changes	future sensory input?

Information bottleneck (IB)



data: $p(x, y)$

free parameter: $\beta > 0$

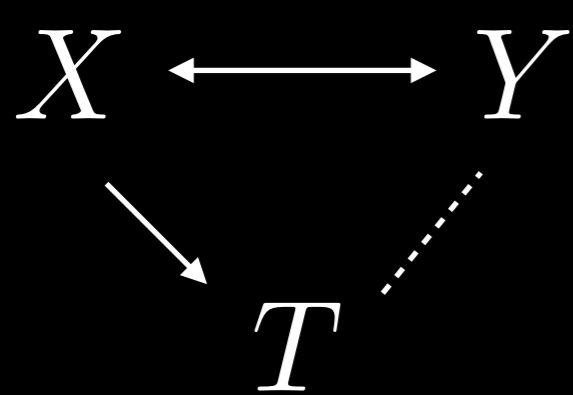
Markov constraint: $T \leftarrow X \longleftrightarrow Y$

$$\min_{q(t|x)} L[q(t|x)] = I(T; X) - \beta I(T; Y)$$

compression

relevance

Information bottleneck (IB)



data: $p(x, y)$

free parameter: $\beta > 0$

Markov constraint: $T \leftarrow X \longleftrightarrow Y$

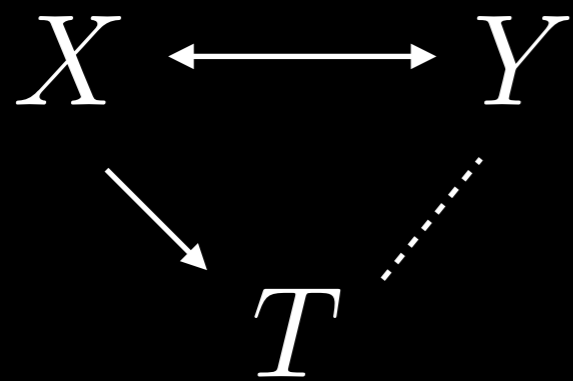
$$\min_{q(t|x)} L[q(t|x)] = I(T; X) - \beta I(T; Y)$$

$$q(t|x) = \frac{q(t)}{Z(x, \beta)} \exp[-\beta D_{KL}[p(y|x) | q(y|t)]]$$

$$q(t) = \sum_x p(x) q(t|x)$$

$$q(y|t) = \frac{1}{q(t)} \sum_x p(y|x) q(t|x) p(x)$$

Information bottleneck (IB)

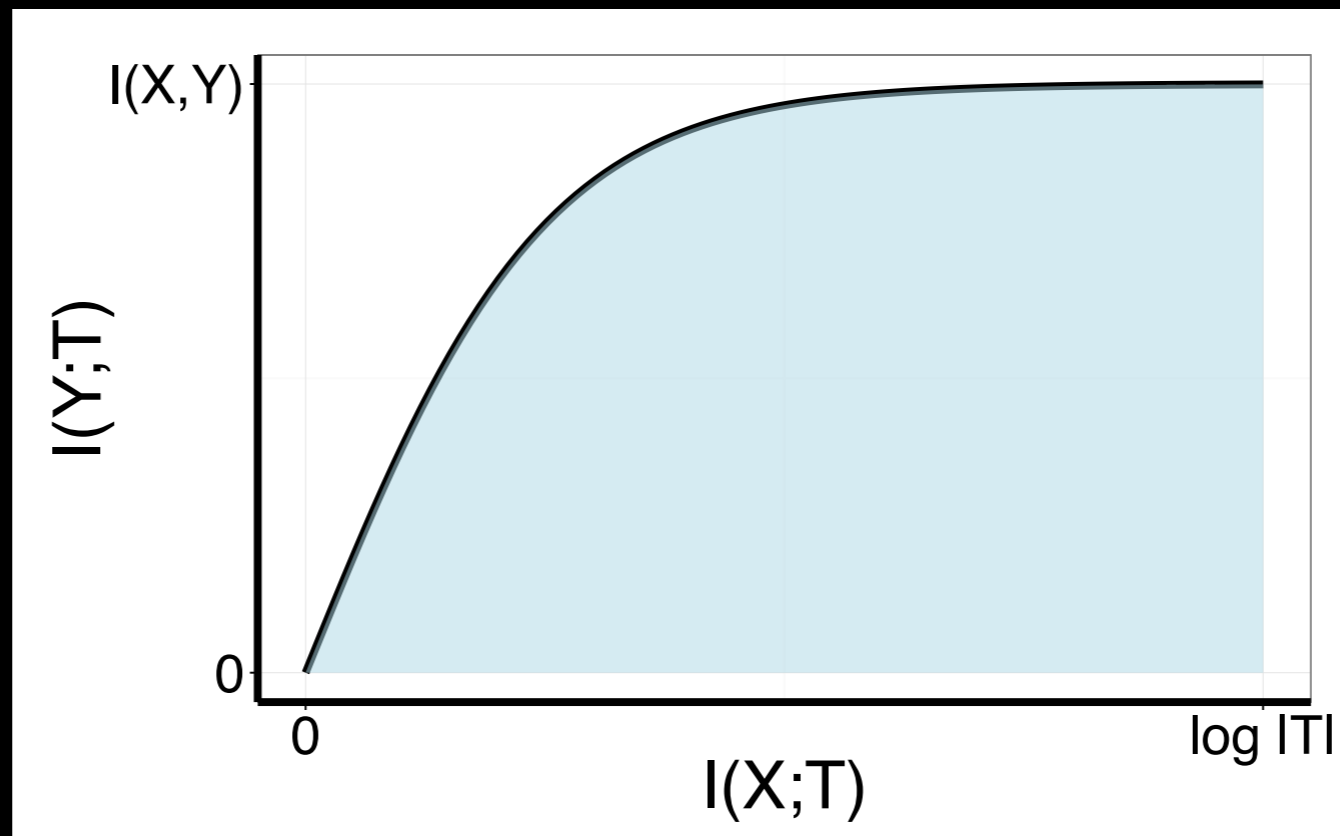


data: $p(x, y)$

free parameter: $\beta > 0$

Markov constraint: $T \leftarrow X \longleftrightarrow Y$

$$\min_{q(t|x)} L[q(t|x)] = I(T; X) - \beta I(T; Y)$$



Measuring compression

$$\min_{q(t|x)} L[q(t|x)] = I(T; X) - \beta I(T; Y)$$

channel coding/
rate distortion theory

$$\min_{q(t|x)} L[q(t|x)] = H(T) - \beta I(T; Y)$$

source coding

$$\begin{aligned} L_{\text{IB}} - L_{\text{DIB}} &= I(X; T) - H(T) \\ &= -H(T | X) \end{aligned}$$

implicit encouragement of stochasticity

A generalized IB

$$L_{\alpha} \equiv H(T) - \alpha H(T | X) - \beta I(Y; T)$$

$$L_{\text{IB}} = L_{\alpha=1}$$

$$L_{\text{DIB}} = L_{\alpha=0}?$$

A generalized IB

$$L_\alpha \equiv H(T) - \alpha H(T | X) - \beta I(Y; T)$$

$$q_\alpha(t | x) \propto \exp \left[\frac{1}{\alpha} (\log q_\alpha(t) - \beta D_{\text{KL}}[p(y | x) | q_\alpha(y | t)]) \right]$$

$$q_{IB}(t | x) = \frac{q(t)}{Z(x, \beta)} \exp[-\beta D_{\text{KL}}[p(y | x) | q(y | t)]]$$

Solving the DIB

$$L_\alpha \equiv H(T) - \alpha H(T | X) - \beta I(Y; T)$$

$$q_\alpha(t | x) \propto \exp \left[\frac{1}{\alpha} (\log q_\alpha(t) - \beta D_{\text{KL}}[p(y | x) | q_\alpha(y | t)]) \right]$$

$$\lim_{\alpha \rightarrow 0} q_\alpha(t | x) = \delta(t - f(x))$$

$$f(x) = \operatorname{argmax}_t (\log q(t) - \beta D_{\text{KL}}[p(y | x) | q(y | t)])$$

IB vs DIB: summary

$$L_{\text{IB}} = I(X; T) - \beta I(Y; T)$$

$$q_{\text{IB}}(t | x) = \frac{q(t)}{Z(x, \beta)} \exp[-\beta D_{\text{KL}}[p(y | x) | q(y | t)]]$$

channel coding with relevance

$$L_{\text{DIB}} = H(T) - \beta I(Y; T)$$

$$q_{\text{DIB}}(t | x) = \delta(t - f(x))$$

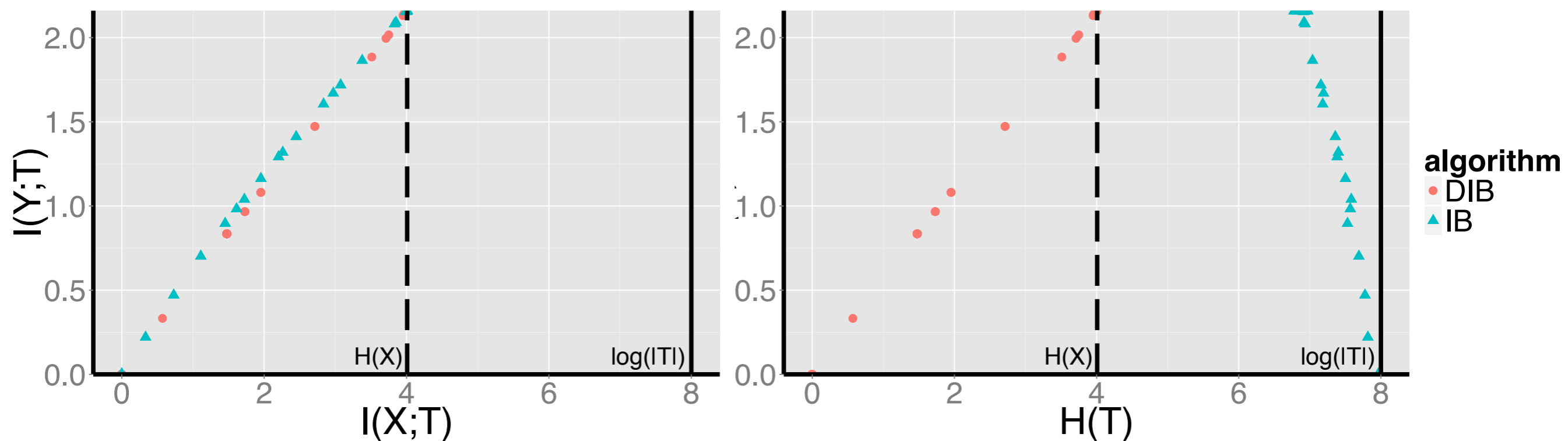
$$f(x) = \underset{t}{\operatorname{argmax}} (\log q(t) - \beta D_{\text{KL}}[p(y | x) | q(y | t)])$$

source coding with relevance

IB vs DIB: experiments

IB plane

DIB plane



$$I(X, T) = H(T) - H(T | X)$$

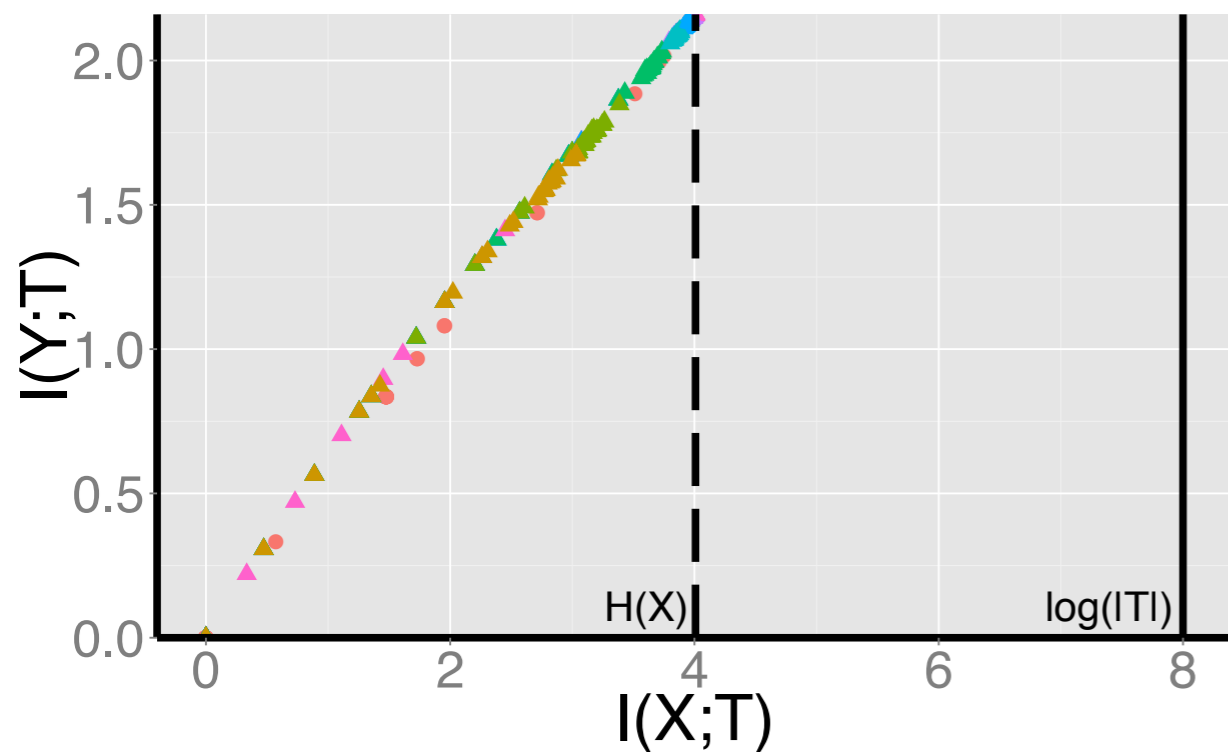
2 ways to get
 $I(X, T) = 0$

$$H(T) = H(T | X) = 0$$

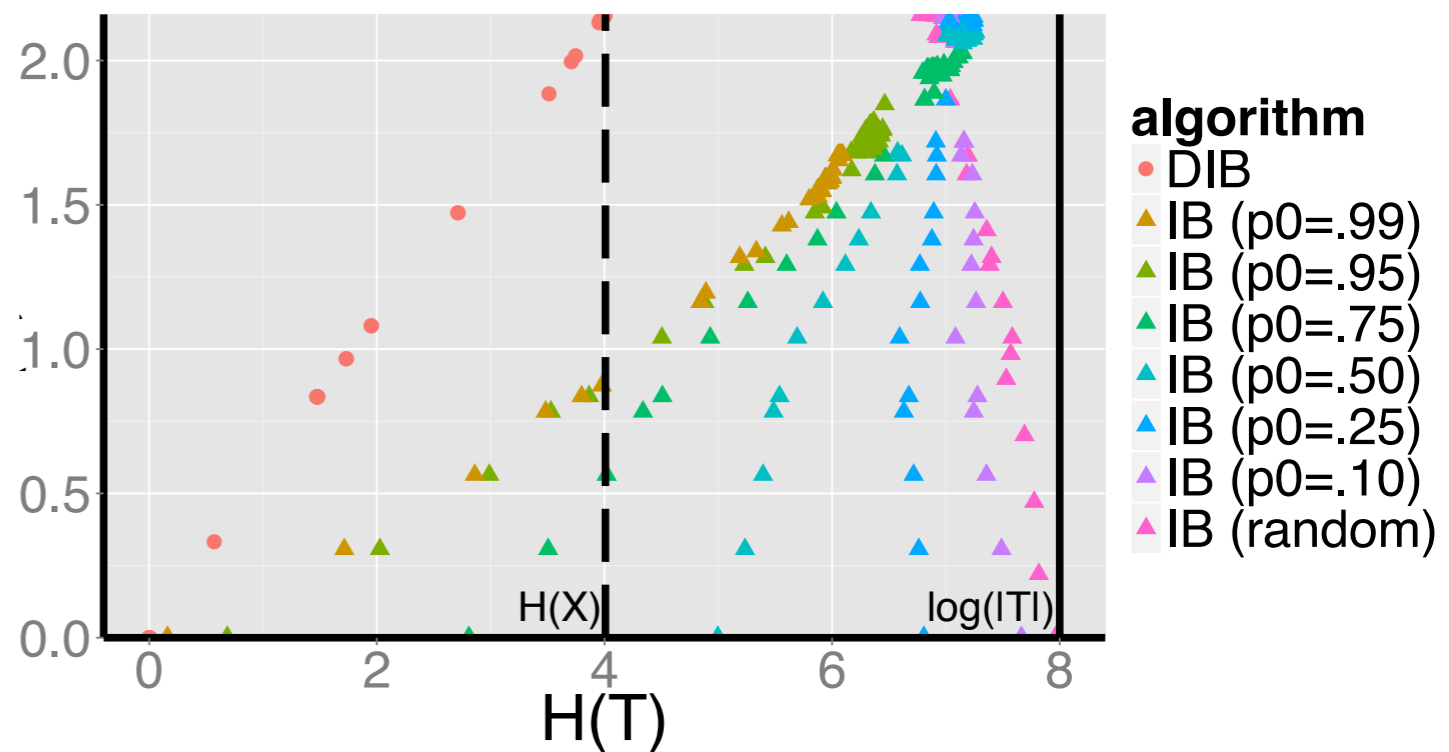
$$H(T) = H(T | X) = \text{const} > 0$$

IB vs DIB: experiments

IB plane

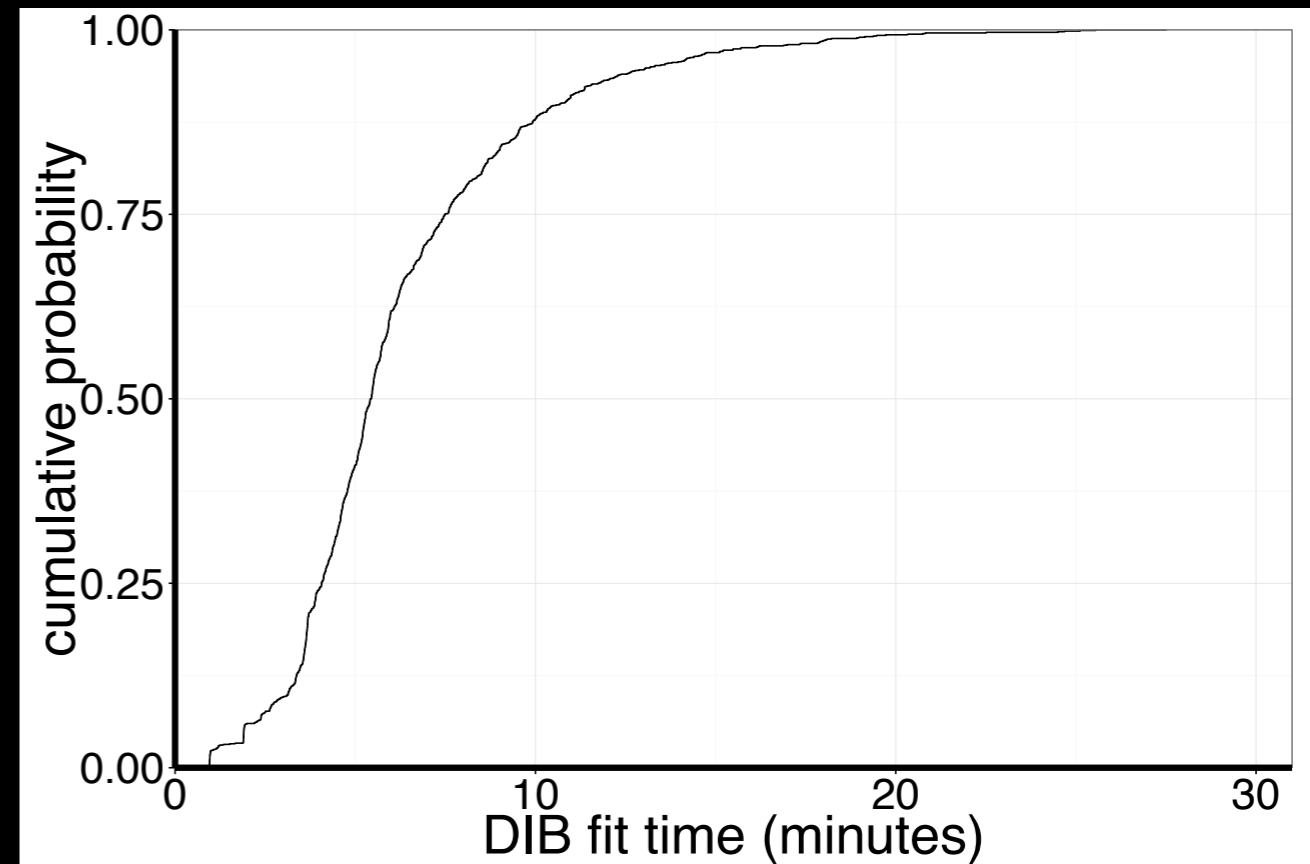
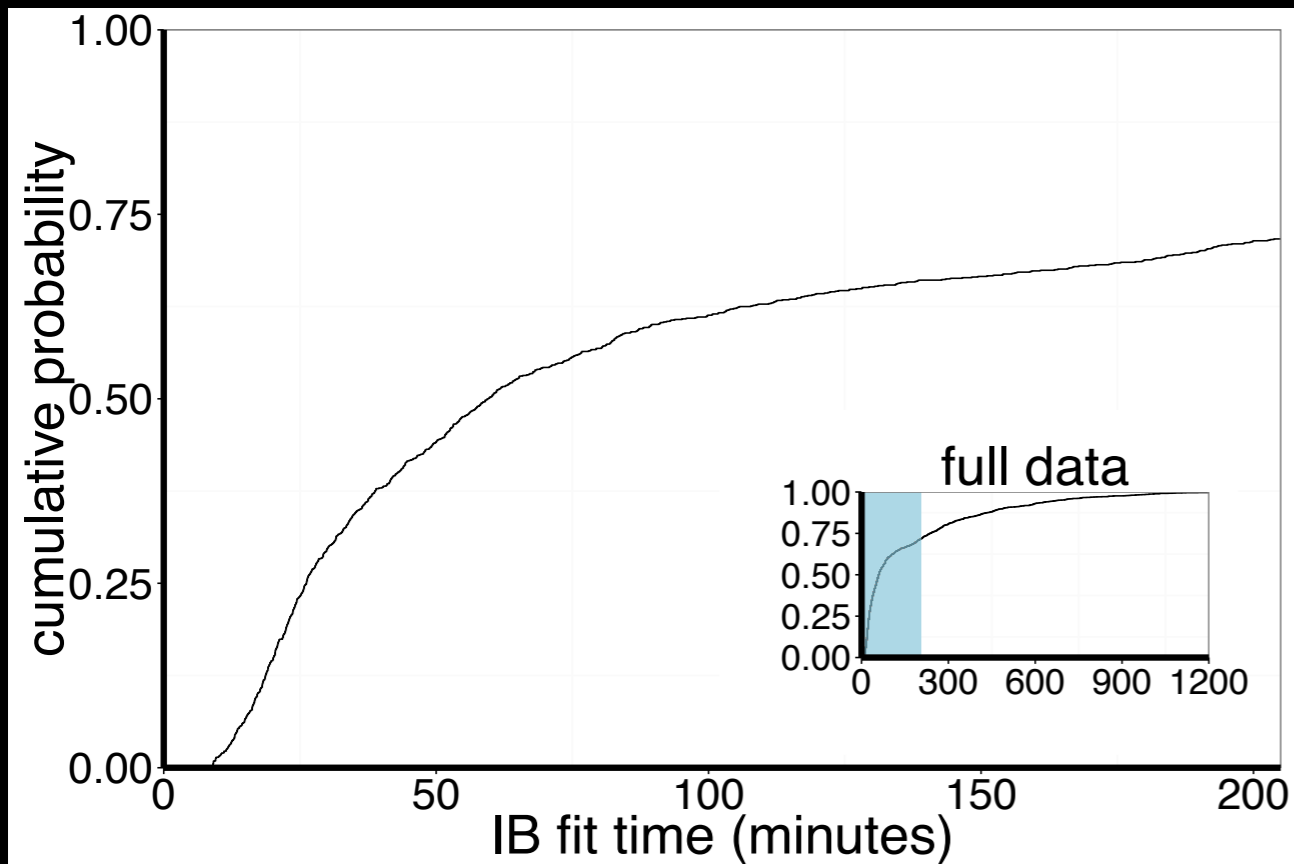


DIB plane

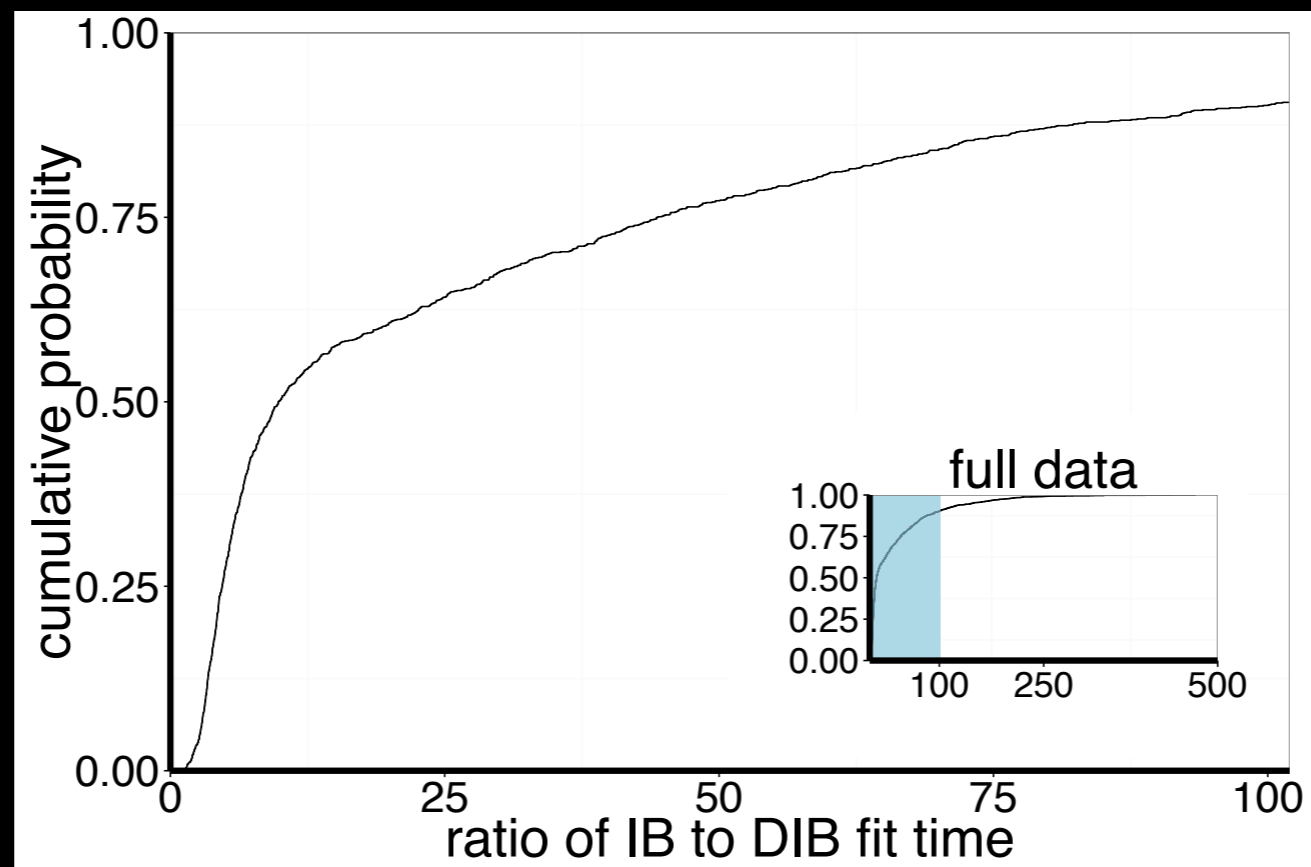


IB and DIB behave similarly in terms of IB cost function...
...but DIB performs far better in terms of its own

IB vs DIB: efficiency



IB vs DIB: efficiency



Summary

- proposed new cost functional for extraction of relevant information based on source coding (rather than channel coding)
- consequence -> deterministic encoder/hard clustering (rather than stochastic/soft)
- IB and DIB exhibit non-trivial differences when fit to data
- DIB fits run 1-2 orders of magnitude faster than IB
- bonus: method to interpolate between IB and DIB

Thanks

- David Schwab & Bill Bialek
- The Krell Institute
- The Hertz Foundation