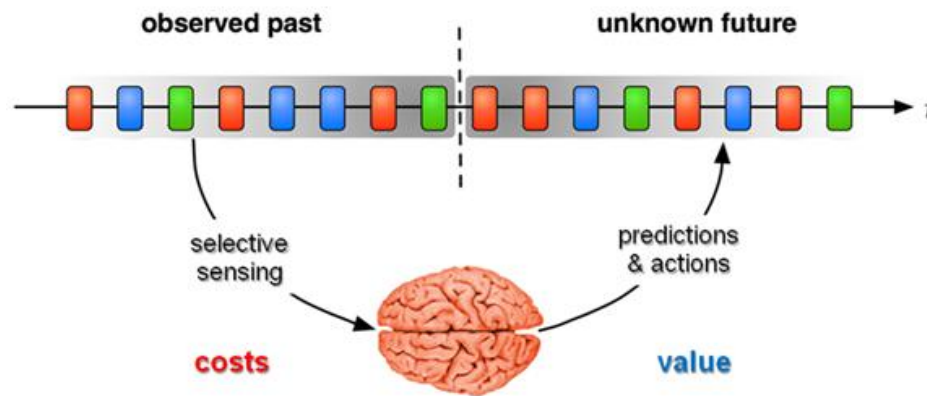


On the balance of Information-flow and Information Discounting in I-RL

Dagstuhl 2011

The brain is a past-future “Information-engine” ...



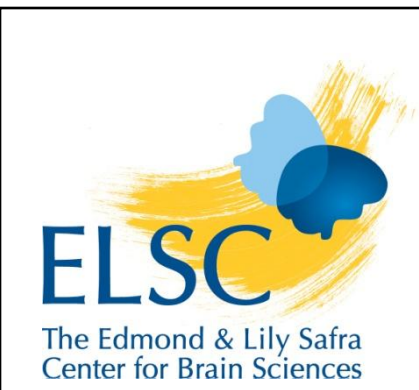
Naftali Tishby



Bill Bialek Daniel Polani Eli Nelken Jonathan Rubin Ohad Shamir

ELSC

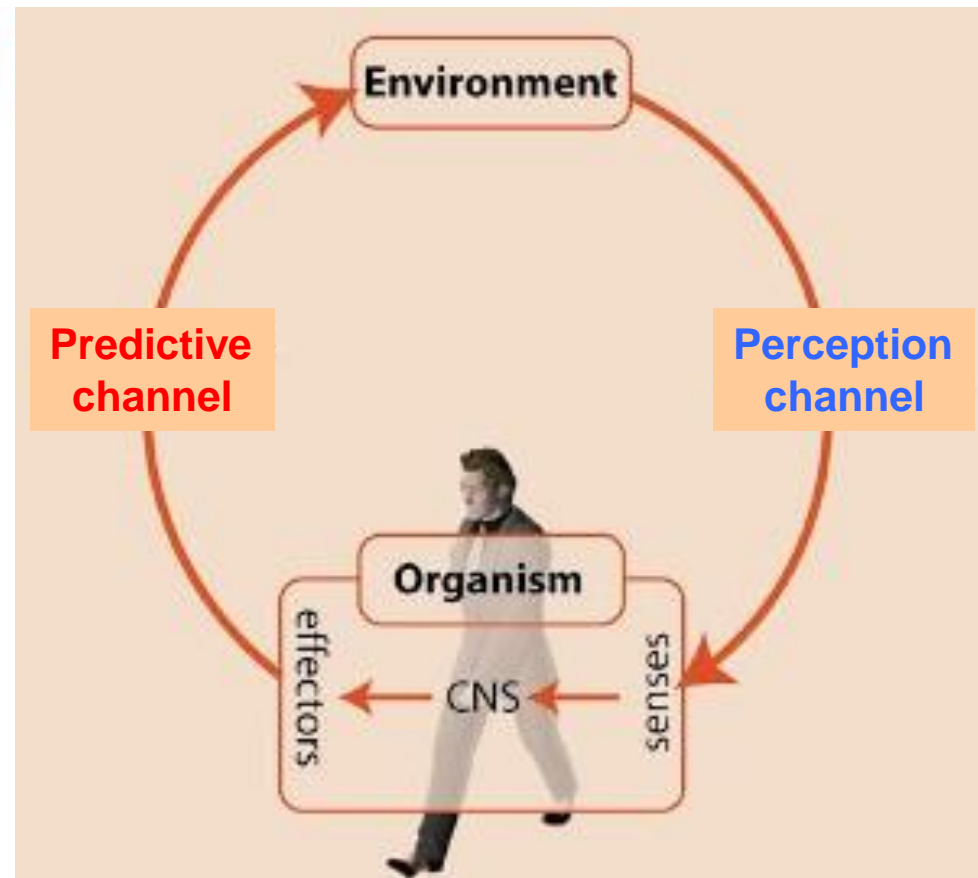
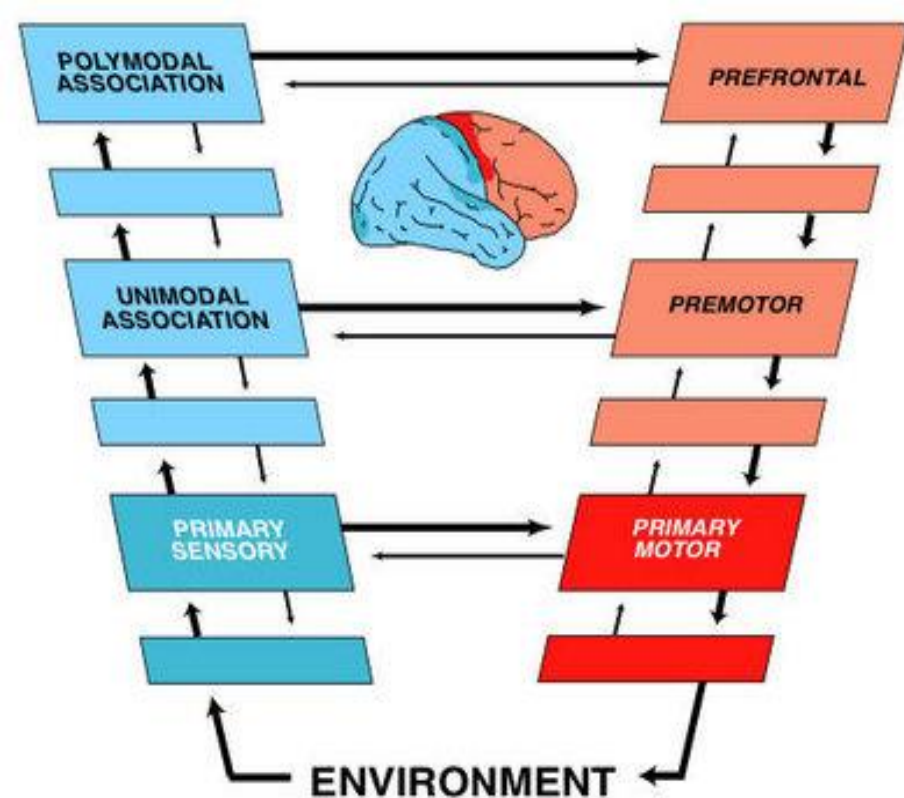
Interdisciplinary Center for Neural Computation
School of Engineering and Computer Science
Hebrew University, Jerusalem, Israel



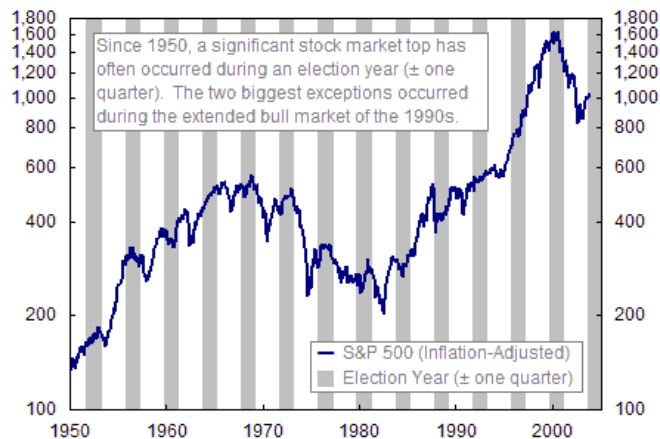
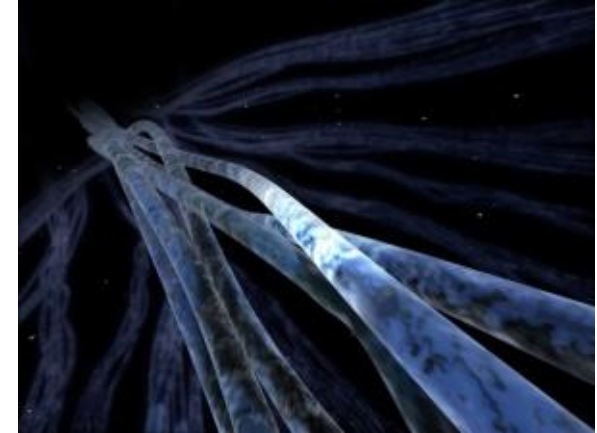
The Perception-Action Cycle

The circular flow of **information** that takes place between the organism and its environment in the course of a sensory-guided sequence of behavior towards a goal.

(JM Fuster)

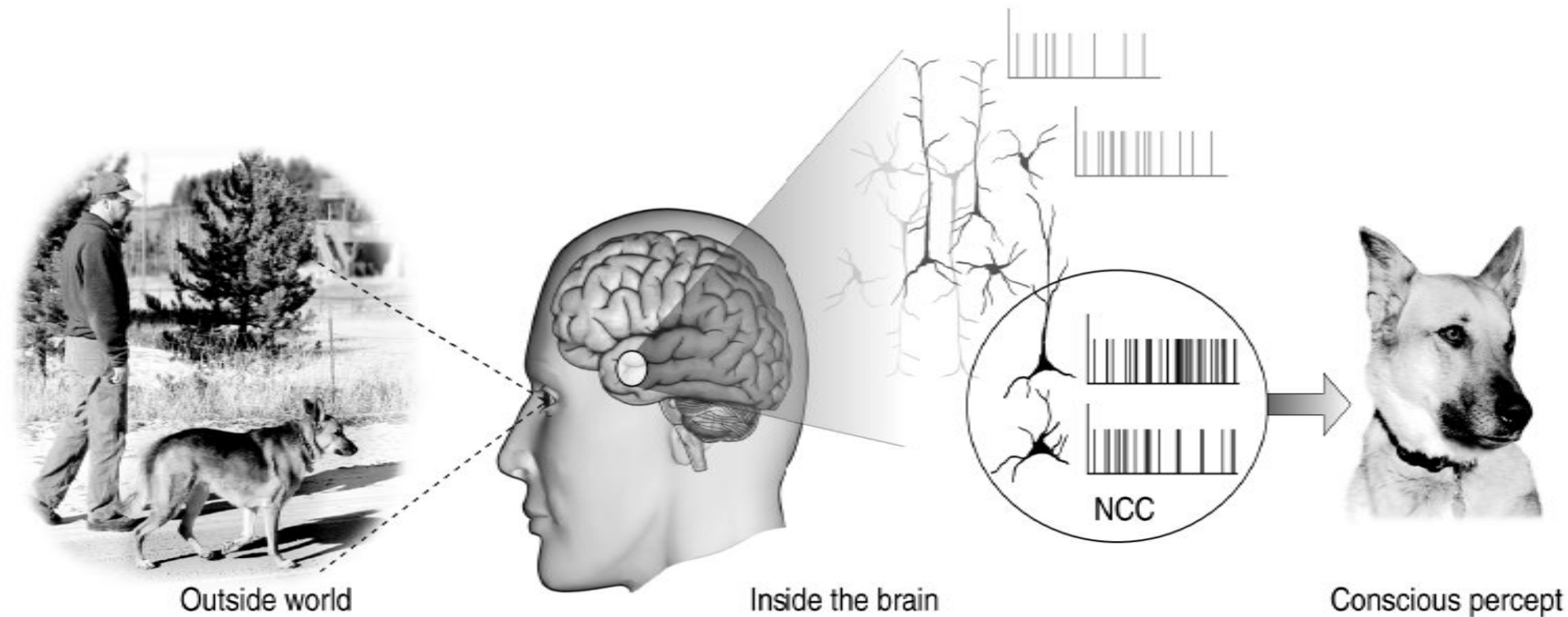


Life is all about making valuable predictions...



The brain's primary task is making valuable predictions

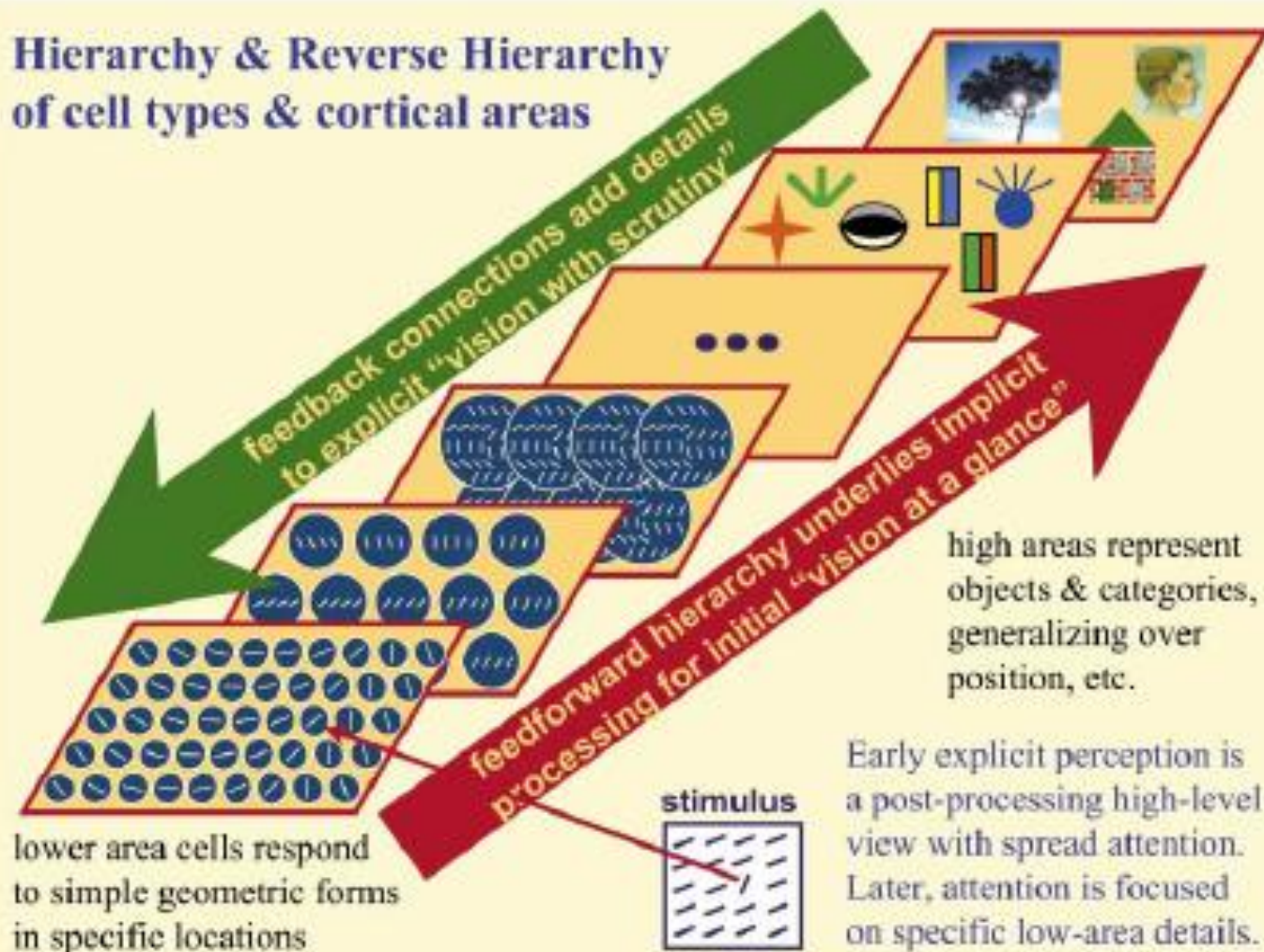
Perception - goal oriented hypotheses generation,
directed by active predictions & useful decisions,
tested by [external or internal] information gathering



From C. Koch,

Hierarchies and reverse hierarchies

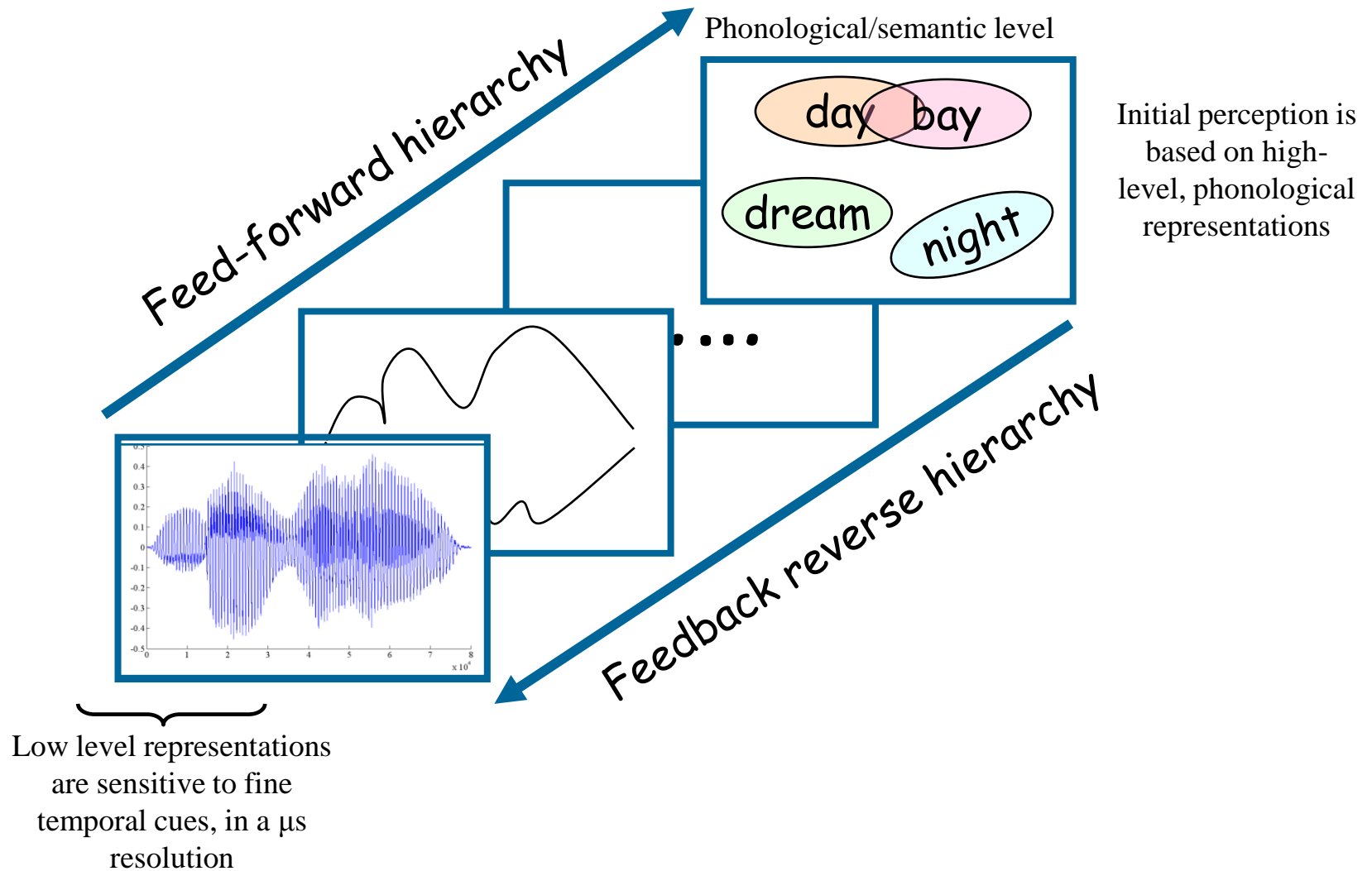
Hierarchy & Reverse Hierarchy of cell types & cortical areas



Tsostos 1990;
Hochstein and
Ahissar 2002

We see what we expect to see





Nelken et al, 2005

[Partially Observed] Markov Decision Processes

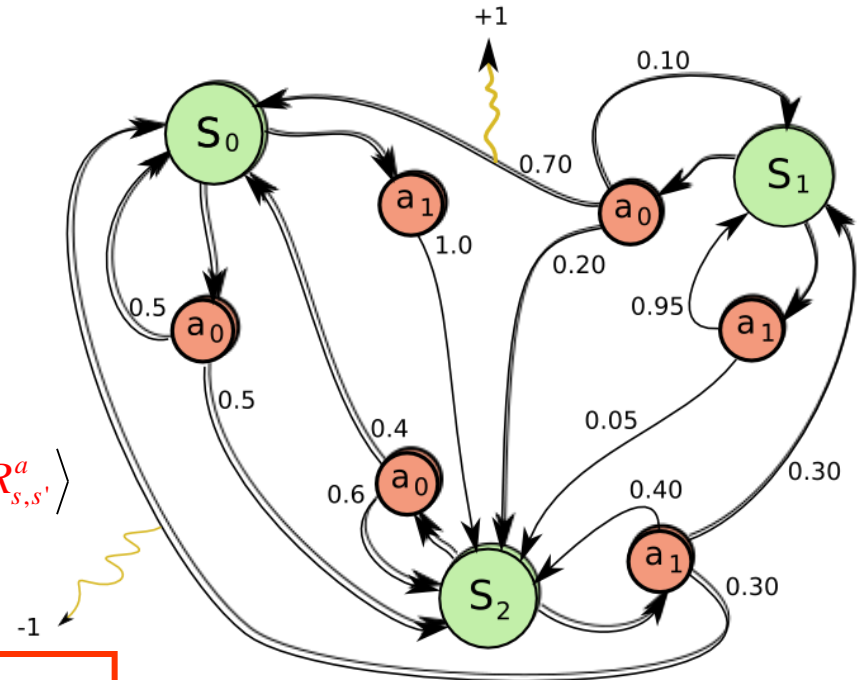
An ergodic MDP : four tuple of

$$\langle s_t \in S, a_t \in A, P_{s,s'}^a = p(s' | s, a), R_{s,s'}^a \rangle$$

states, actions, ergodic transition probabilities and rewards.

An ergodic POMDP also has stochastic observations at each state (as in HMM)

$$\langle s_t \in S, o_t \in O, \sigma(o | s), a_t \in A, P_{s,s'}^a = p(s' | s, a), R_{s,s'}^a \rangle$$



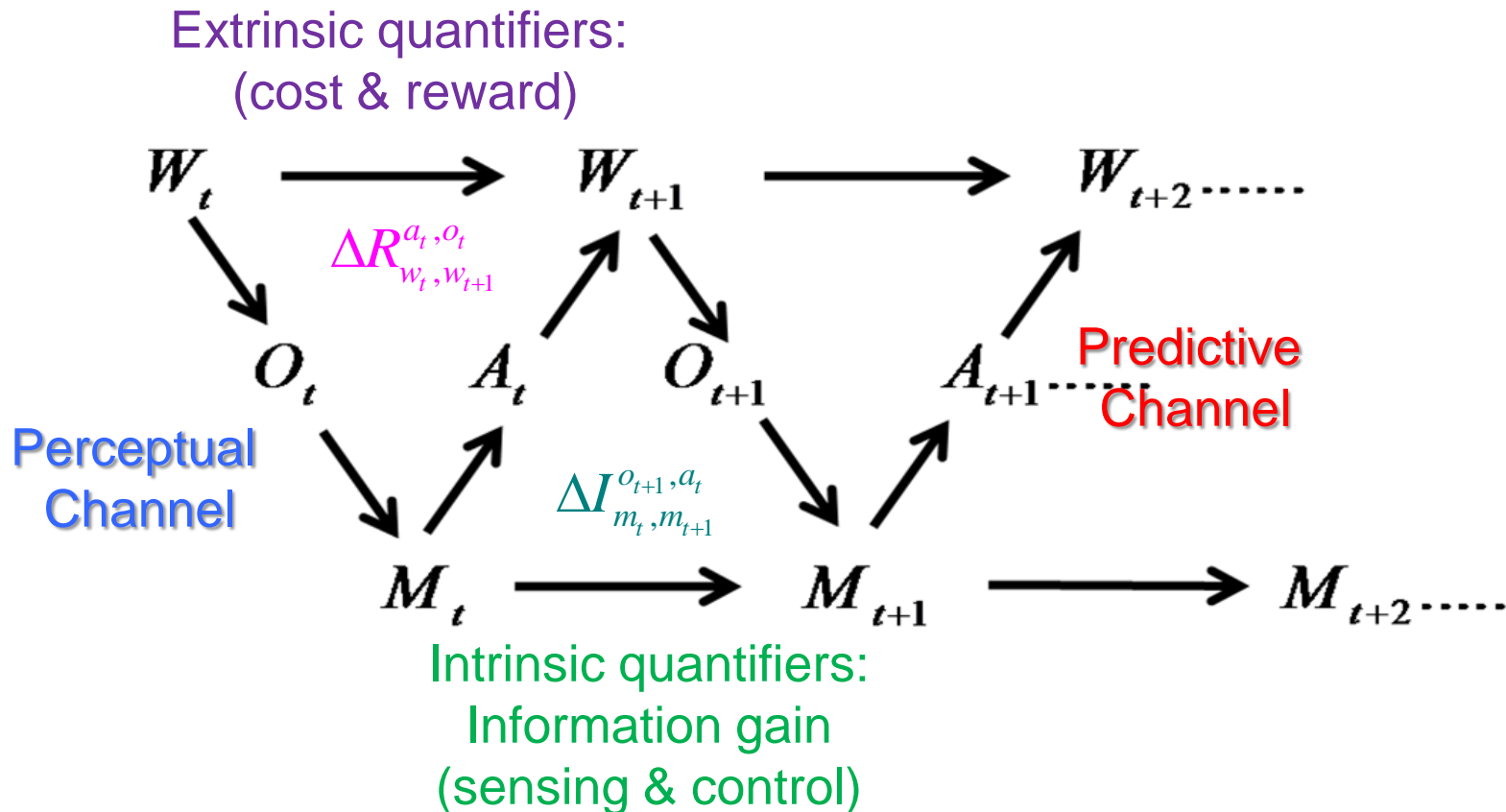
The planning problem :

Find the optimal [?] policy, $\pi(a | s)$, that maximizes

expected future reward: $E_{p(a_t, s_{t+1}, a_{t+1}, s_{t+2} \dots | s_t)} \left[\sum_{j=t}^{\infty} R_{s_j, s_{j+1}}^{a_j} \right]$

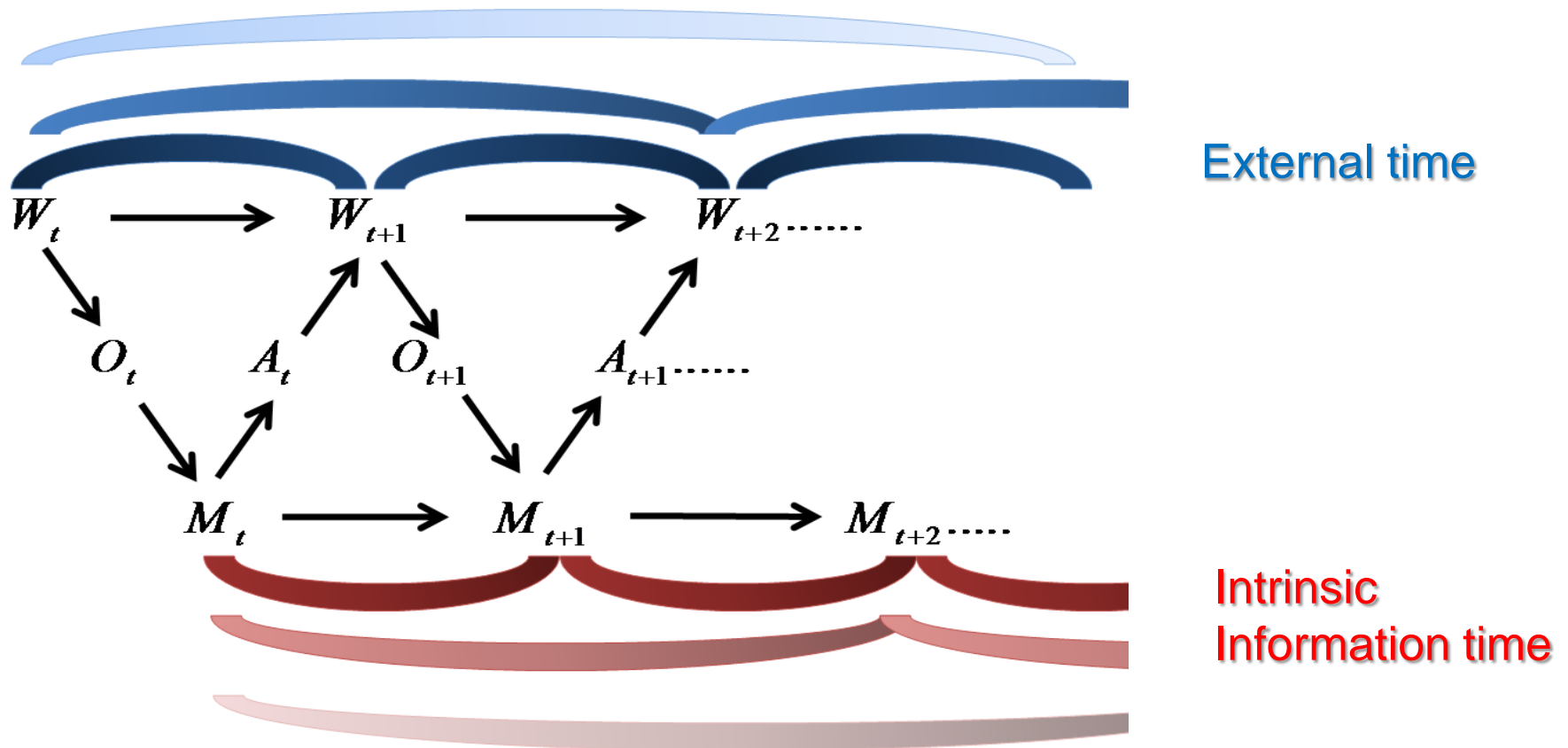
subject to various constraints.

Graphical model for the perception-action-cycle



Both future extrinsic reward (Value) and intrinsic (Information to-go) are optimized together using Bellman-like equations w.r.t. to both channels.

POMDPs, HMMs and Predictive Information



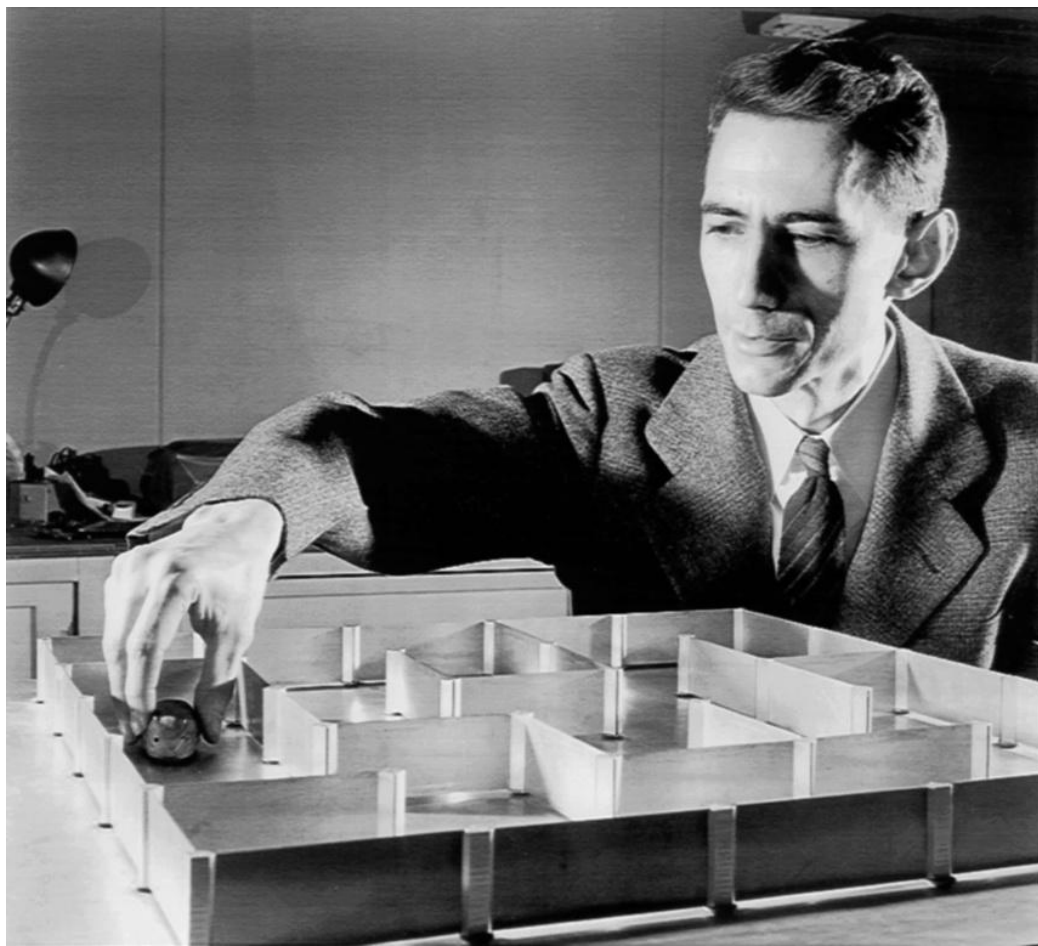
Top-down perceptual & planning hierarchy improves the information gain and the predictive horizon.

Correct discounting of past and future is a natural consequence.

Bellman meets Shannon



Richard Ernest Bellman
(August 26, 1920 – March 19, 1984)



Claude Elwood Shannon
(April 30, 1916 – February 24, 2001)

Decision-sequences and information

Proposition 1 :

state/decision sequences behave like codes:

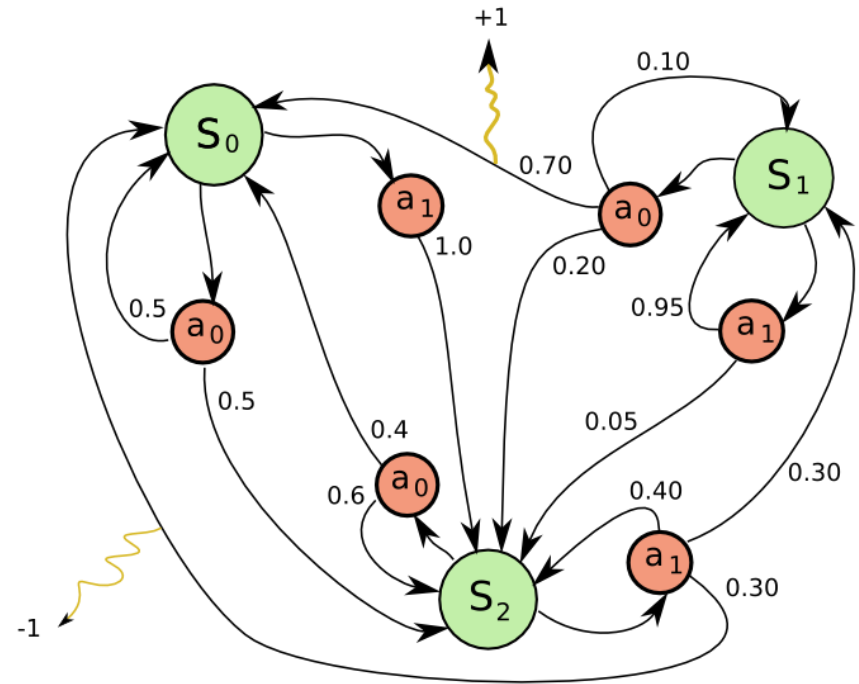
- they can be concatenated to longer sequences
- sequence lengths obey a **Kraft inequality**
- number/**complexity of** (binary) **decisions** along a trajectory, is **lower-bounded** by the **decision Entropy** (as in source coding).

Proposition 2 :

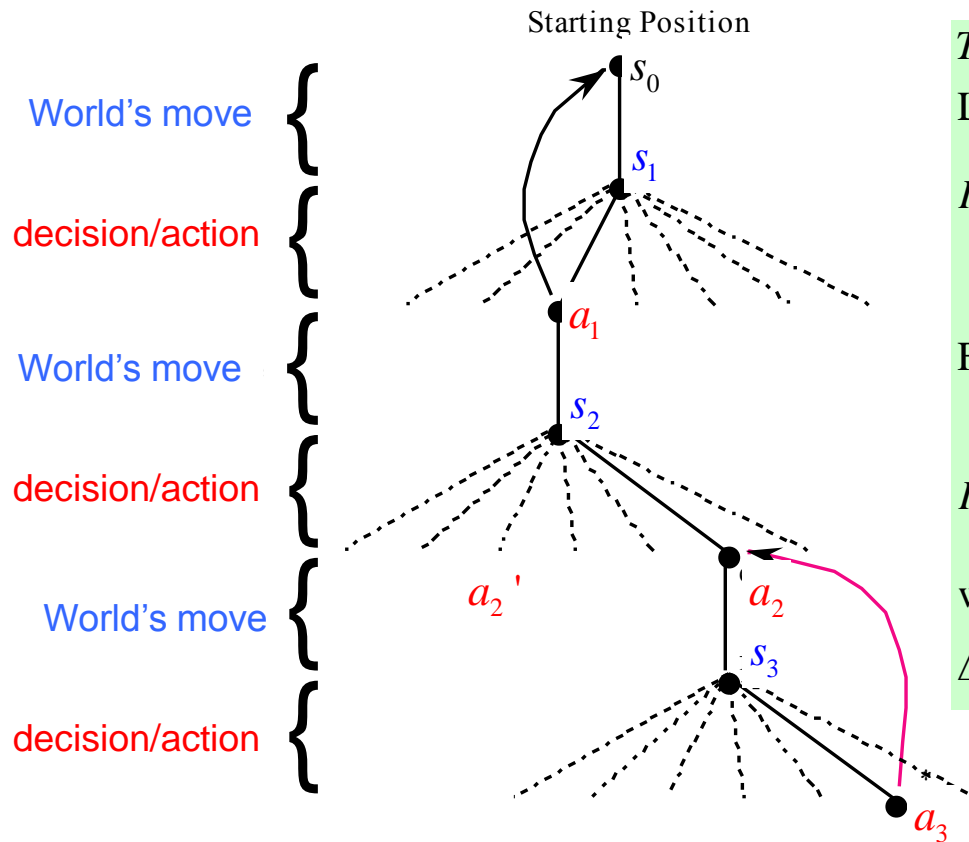
policies are like code nodes
subsequences of optimal decision sequences
are also optimal - **Bellman optimality condition**.

Proposition 3 :

Entropy (and relative-entropy, KL-divergences)
are **the only** (up to units) functionals of the state-decision
sequence pdf that obey an **additive Bellman-like equation**
on the simplex.



Decision/action sequences and information



s - the state before our action

s' - the state after our action

Theorem :

Let G denote our target (relevant) variable

$$I(s; G) = E_g \log \frac{p(g | s)}{p(g)} = \sum_{g \in G} p(g, s) \log \frac{p(g | s)}{p(g)}$$

- the **specific** mutual information at state s on G .

For an MDP $\langle S, A, P, \pi \rangle$ the following recursion holds:

$$I^\pi(s_t; G) = \sum_{a_t \in A} \pi(a_t | s_t) \sum_{s_{t+1}} P^{a_t}_{s_t, s_{t+1}} [\Delta I^{a_t}_{s_t, s_{t+1}} + I^\pi(s_{t+1}; G)]$$

with

$$\Delta I^{a_t}_{s_t, s_{t+1}} \triangleq -I(G; s_{t+1} | s_t, a_t)$$

Proof idea: Recursive Information-chain rules...

The future actions-states in MDP sequence: $A_t, S_{t+1}A_{t+1}, S_{t+2}A_{t+2}, \dots, S_{t+n}A_{t+n}, \dots$

Let: $\mathcal{J}^\pi(s_t, a_t) \triangleq E_{p(s_{t+1}, a_{t+1}, s_{t+2}, a_{t+2}, \dots | s_t, a_t)} \log \frac{p(s_{t+1} a_{t+1} s_{t+2} a_{t+2}, \dots | s_t, a_t)}{p(s_{t+1}) \pi(a_{t+1}) p(s_{t+2}) \pi(a_{t+2}), \dots}$

then:

$$\mathcal{J}^\pi(s_t, a_t) = E_{p(s_{t+1}, a_{t+1} | s_t, a_t)} \left[\log \frac{p(s_{t+1} | s_t, a_t)}{p(s_{t+1})} + \log \frac{\pi(a_{t+1} | s_{t+1})}{\pi(a_{t+1})} + \mathcal{J}^\pi(s_{t+1}, a_{t+1}) \right]$$

With: $p(s_{t+1}, a_{t+1} | s_t, a_t) = p(s_{t+1} | s_t, a_t) \pi(a_{t+1} | s_{t+1})$

$$\text{and: } \Delta I_{s_t, s_{t+1}}^{a_t} \triangleq \log \frac{p(s_{t+1} | s_t, a_t)}{p(s_{t+1})} + \log \frac{\pi(a_{t+1} | s_{t+1})}{\pi(a_{t+1})}$$

□

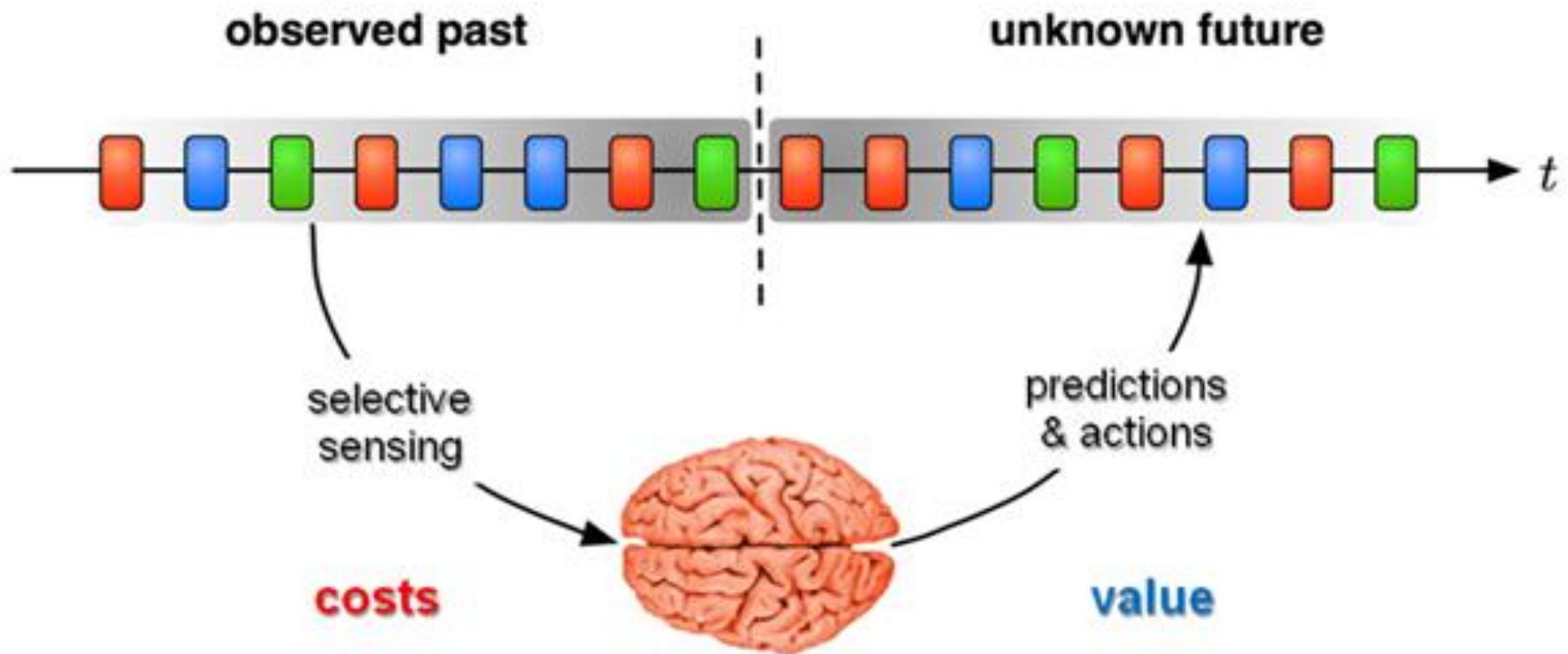
$\log \frac{p(s_{t+1} | s_t, a_t)}{p(s_{t+1})}$ - sensory information gain (perception capacity)

$\log \frac{\pi(a_{t+1} | s_{t+1})}{\pi(a_{t+1})}$ - **required** predictive capacity (control complexity)

Predictive Information:
Information in the past about the future

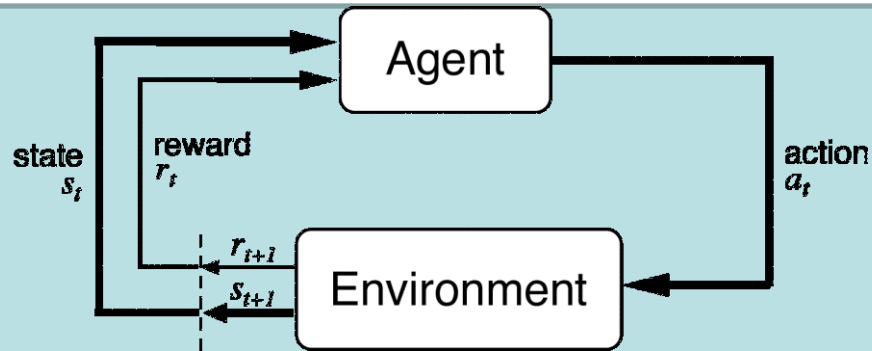
The essence of the cycle

The brain is a past-future “Information-engine” ...



**How much information is needed
for valuable behavior?**

Value (extrinsic) and Information (intrinsic)...



Agent and environment interact at time steps: $t = 0, 1, 2, \dots$

Agent observes TRUE state at step t : $s_t \in S$

produces action at step t : $a_t \in A(s_t)$ with $\pi(a | s)$

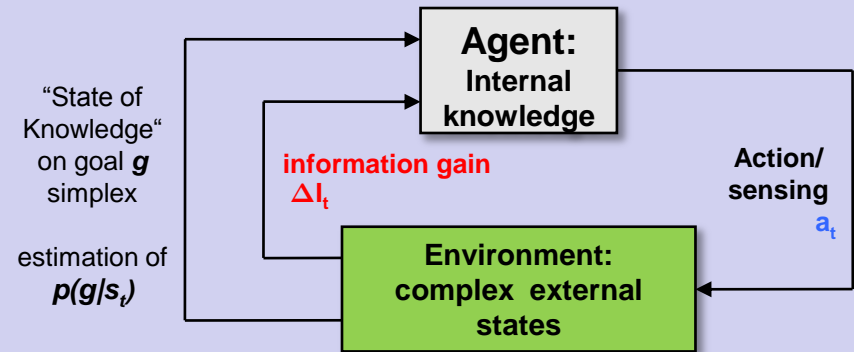
gets resulting reward: $r_{t+1} \in \mathcal{R}$

resulting next state: s_{t+1} with $p(s_{t+1} | s_t, a_t) = P_{s,s'}^a$

Bellman equation for value V^π :

$$V^\pi(s) = \sum_a \pi(a | s) \sum_{s'} P_{ss'}^a [R_{ss'}^a + \gamma V^\pi(s')]$$

solved for $V^\pi(s)$ by DP given $P_{ss'}^a, R_{ss'}^a$ and π



Agent has a goal variable $g \in G$

interacts with environment at time steps: $t = 0, 1, 2, \dots$

estimates/infer an internal state: $\hat{s}_t \in \hat{S}$,

characterized by $p(\hat{s} | s)$, $p(g | \hat{s})$

produces action at step t : $a_t \in A(s_t)$ with $\pi(a | \hat{s})$

get/estimate information gain: $\Delta I_{s,s'}^a \in \mathcal{R}$

resulting world next state: s_{t+1} with $P_{s,s'}^a$

Bellman equation for I^π :

$$I^\pi(\hat{s}; g) = \sum_a \pi(a | \hat{s}) \sum_{s'} P_{ss'}^a [\Delta I_{ss'}^a + I^\pi(\hat{s}'; g)]$$

solved for I^π using DP and prob. inference

Combining (future) Value and Information

In cases where information is free, we can maximize value irrespective of its information cost.

In general, however, we want

- (1) **to reduce decision complexity** (get home in the simplest way)
- (2) **maximize the environment information gain** (e.g. with the coins)
- (3) **increase robustness to model fluctuations**

All three can be obtained by combining the Information and Value equations.

Trading Value and (future) Information

$$\mathcal{J}^\pi(s_t, a_t) \triangleq E_{p(s_{t+1}, a_{t+1}, s_{t+2}, a_{t+2}, \dots | s_t, a_t)} \log \frac{p(s_{t+1} a_{t+1} s_{t+1} a_{t+1}, \dots | s_t, a_t)}{p(s_{t+1}) \pi(a_{t+1}) p(s_{t+2}) \pi(a_{t+2}), \dots} \text{ then:}$$

$$\mathcal{J}^\pi(s_t, a_t) = E_{p(s_{t+1}, a_{t+1} | s_t, a_t)} \left[\log \frac{p(s_{t+1} | s_t, a_t)}{p(s_{t+1})} + \log \frac{\pi(a_{t+1} | s_{t+1})}{\pi(a_{t+1})} + \mathcal{J}^\pi(s_{t+1}, a_{t+1}) \right]$$

$$\text{With: } p(s_{t+1}, a_{t+1} | s_t, a_t) = p(s_{t+1} | s_t, a_t) \pi(a_{t+1} | s_{t+1})$$

$$\text{We want: } \arg \min_{\pi} \mathcal{J}^\pi(s_t, a_t) - \beta Q^\pi(s_t, a_t) \triangleq \arg \min_{\pi} F^\pi(s_t, a_t, \beta)$$

$$\text{with: } Q^\pi(s_t, a_t) = \sum_{s_{t+1}} p(s_{t+1} | s_t, a_t) [R_{s_t s_{t+1}}^{a_t} + V^\pi(s_{t+1})]$$

$$\mathcal{J}^\pi(s_t, a_t) - \beta Q^\pi(s_t, a_t) =$$

$$= E_{p(s_{t+1} | s_t, a_t) \pi(a_{t+1} | s_{t+1})} \left[\log \frac{p(s_{t+1} | s_t, a_t)}{p(s_{t+1})} + \log \frac{\pi(a_{t+1} | s_{t+1})}{\pi(a_{t+1})} - \beta R_{s_t s_{t+1}}^{a_t} + \mathcal{J}^\pi(s_{t+1}, a_{t+1}) - \beta V^\pi(s_{t+1}) \right]$$

$$= E_{p(s_{t+1} | s_t, a_t) \pi(a_{t+1} | s_{t+1})} \left[\log \frac{p(s_{t+1} | s_t, a_t)}{p(s_{t+1})} + \log \frac{\pi(a_{t+1} | s_{t+1})}{\pi(a_{t+1})} - \beta R_{s_t s_{t+1}}^{a_t} + \mathcal{J}^\pi(s_{t+1}, a_{t+1}) - \beta Q^\pi(s_{t+1}, a_{t+1}) \right]$$

or

$$F^\pi(s_t, a_t, \beta) = E_{p(s_{t+1} | s_t, a_t)} \left[\log \frac{p(s_{t+1} | s_t, a_t)}{p(s_{t+1})} - \beta R_{s_t s_{t+1}}^{a_t} + E_{\pi(a_{t+1} | s_{t+1})} \left[\log \frac{\pi(a_{t+1} | s_{t+1})}{\pi(a_{t+1})} + F^\pi(s_{t+1}, a_{t+1}, \beta) \right] \right]$$

Information bounded RL

define : $q(s' | s, a) = q_{s,s'}^a = \frac{p(s')}{Z(\beta, s, a)} \exp(\beta R_{s,s'}^a)$

the "optimal" (reward as **sufficient statistic**) transition probabilities,
and $p(s')$ the state prior,

$Z(\beta, s, a) = \sum_{s'} p(s') \exp(\beta R_{s,s'}^a)$, is the local partition function.

Then the state-action free energy, $F^\pi(s_t, a_t, \beta) \triangleq \mathcal{J}^\pi(s_t, a_t) - \beta Q^\pi(s_t, a_t)$, Bellman equation is:

$$F^\pi(s_t, a_t, \beta) = \sum_{s_{t+1}} p(s_{t+1} | s_t, a_t) \sum_{a_{t+1}} \pi(a_{t+1} | s_{t+1}) \left[\log \frac{p(s_{t+1} | s_t, a_t)}{q(s_{t+1} | s_t, a_t)} + \log \frac{\pi(a_{t+1} | s_{t+1})}{\pi(a_{t+1})} \right. \\ \left. + F^\pi(s_{t+1}, a_{t+1}, \beta) - \log Z(\beta, s_t, a_t) \right]$$

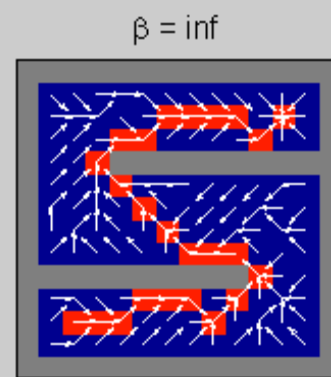
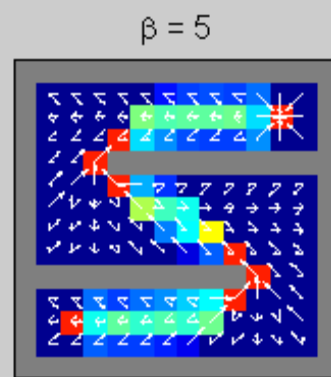
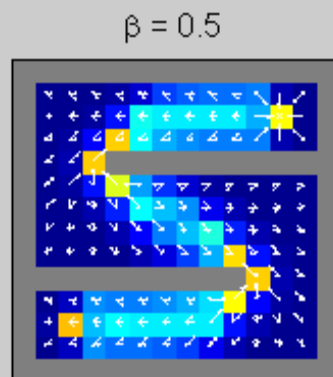
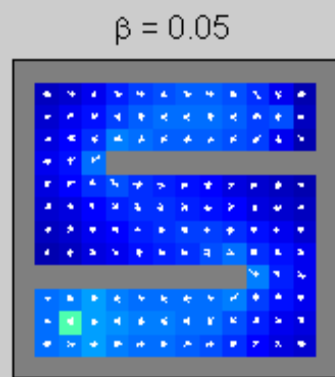
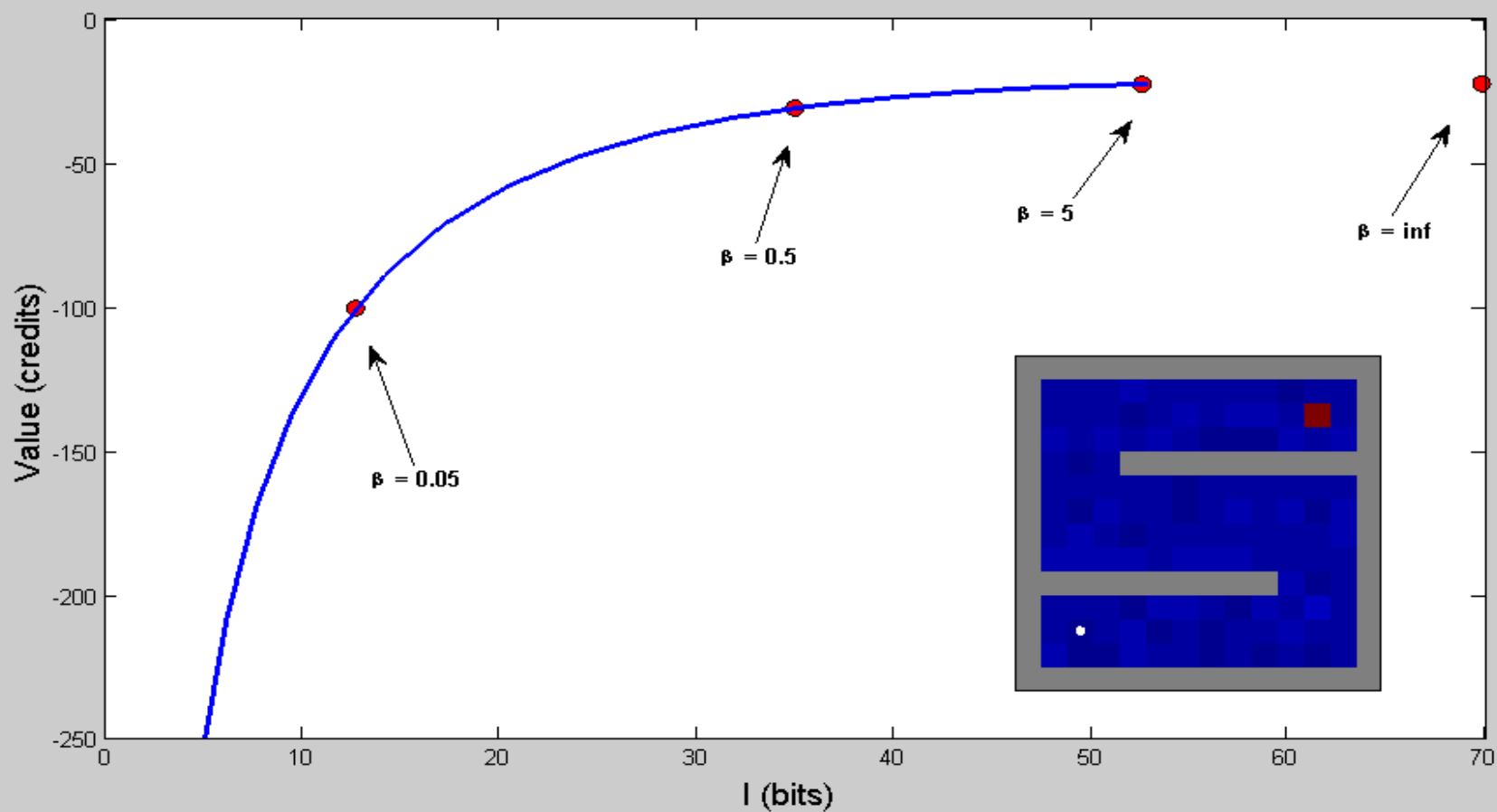
the desired optimal policy is (somewhat surprisingly):

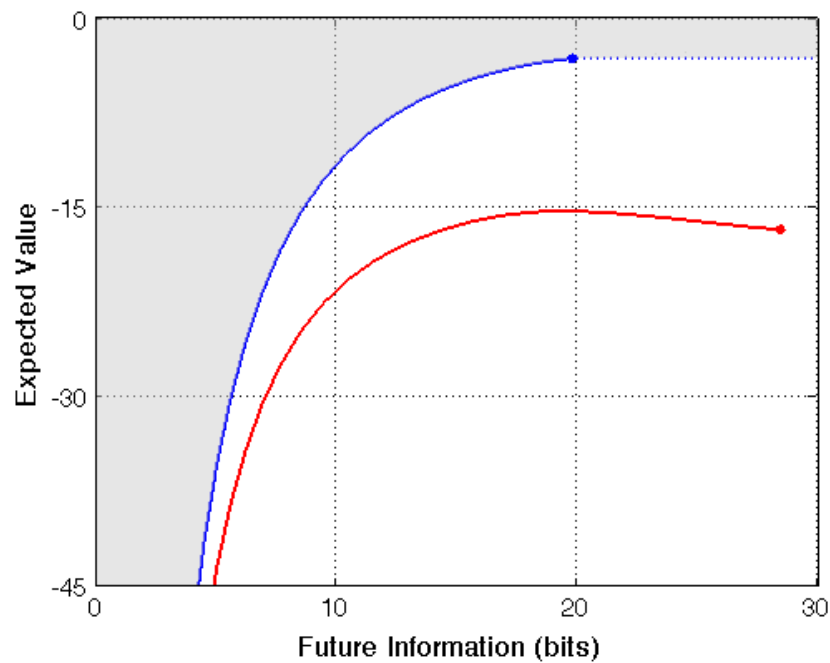
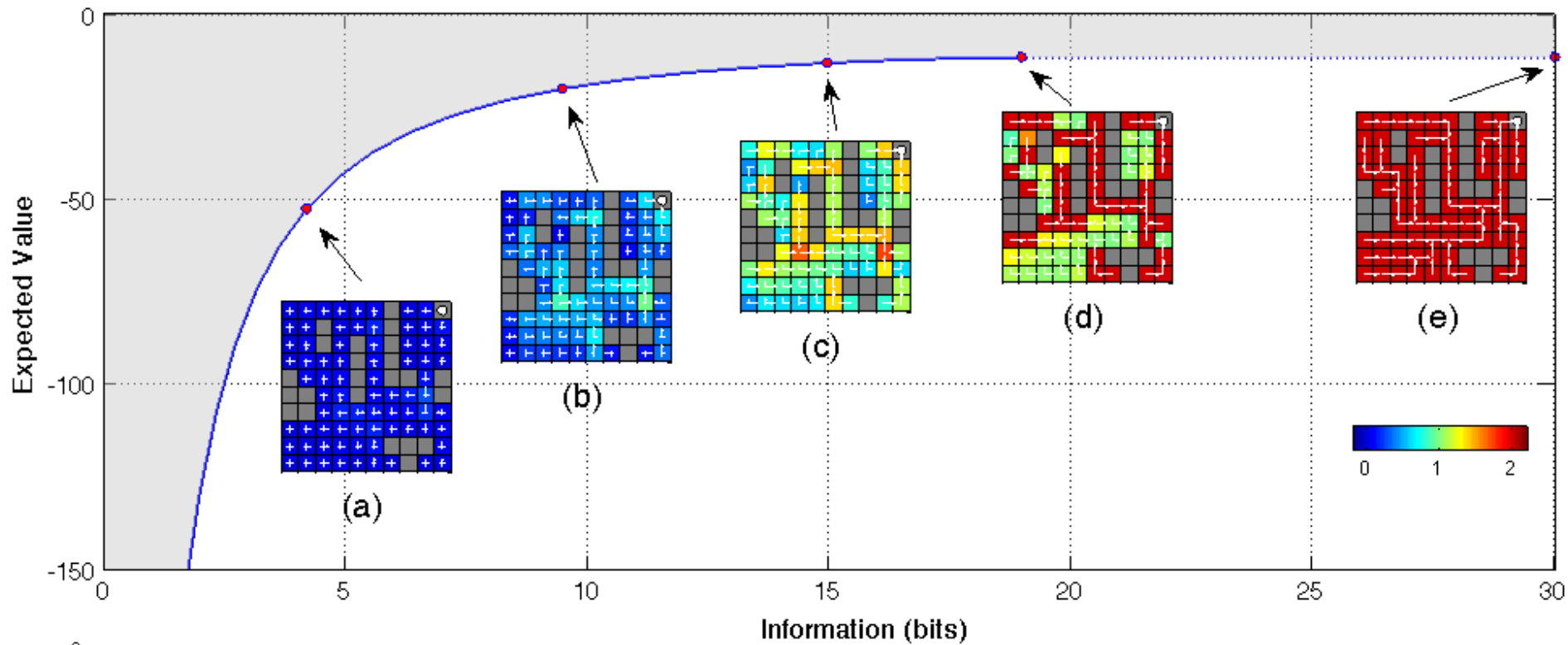
$$\pi(a | s) = \frac{\pi(a)}{Z(s, \beta)} \exp(F^\pi(s, a, \beta))$$

$$Z(s, \beta) = \sum_a \pi(a) \exp(F^\pi(s, a, \beta))$$

$$\pi(a) = \sum_s \pi(a | s) p(s)$$

These 3 equations should be iterated till convergence for every state (like Blahut Arimoto).



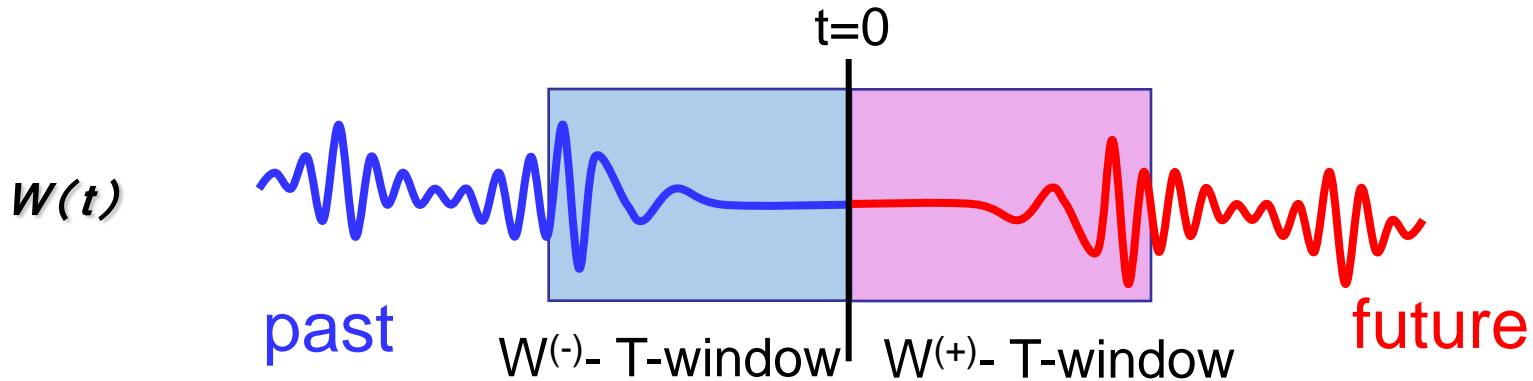


- The optimal tradeoff between Value and Future Information.
- Lower information \rightarrow softer policy
- The optimal – out of training – value is obtained at a finite β

**How to discount
Information (intrinsic) Rewards?**

Predictive Information: The Capacity of the Future-Past Channel

(with Bialek and Nemenman, 2001)



- Estimate $P^T(W^{(-)}, W^{(+)})$: T -past-future distribution



$$I_{pred}[T] = \left\langle \log \frac{p(W_{future}^T | W_{past}^T)}{p(W_{future}^T)} \right\rangle_{p(W_{past}, W_{future})}$$

Logarithmic growth for finite dimensional processes

- Finite parameter processes
(e.g. Markov chains)

$$I_{pred}(T \rightarrow \infty) \approx \frac{\dim(\theta)}{2} \log T$$

- Similar to stochastic complexity (MDL)

Power law growth

- Fast growth is a signature of infinite dimensional processes (e.g. speech)

$$I_{pred}(T \rightarrow \infty) \approx T^{\alpha} \quad \alpha < 1$$

- Power laws appear in cases where the interactions/correlations have long range.

Information gain discounting

- Information gains should accumulate sub-linearly
 \Rightarrow discounting by a factor $\gamma(t) = t^{-\eta}$, $0 \leq \eta \leq 1$
- $\eta = 1 \Rightarrow$ Logarithmic predictive information
fixed intrinsic/extrinsic rewards ratio:
 \Rightarrow exponential discounting of external rewards
 \Rightarrow logarithmic "Intrinsic Time Fisheye"
- $\eta < 1 \Rightarrow$ Power-law predictive information
 \Rightarrow hyperbolic discounting of external rewards
 \Rightarrow hyperbolic intrinsic "Time Fisheye"



Time Fisheye