

Appendix

A Simplified generative model

(i) **Emitters independence** A first assumption is that two transition processes are independent, i.e. $\Pr(Z, U) = \Pr(Z)\Pr(U)$ so that:

$$(\mathbf{z}, \mathbf{u}) \sim \Pr(Z, U | \mathbf{a}, \mathbf{z}_0, \mathbf{u}_0) = \Pr(Z | \mathbf{a}, \mathbf{z}_0, \mathbf{u}_0) \Pr(U | \mathbf{a}, \mathbf{z}_0, \mathbf{u}_0)$$

(ii) **End-effector control** An additional assumption is that the controlled transition process is relatively “fast” in comparison with the uncontrolled one (for, e.g, saccades can be realized in a 100-200 ms interval). In consequence we assimilate the motor command \mathbf{a} with a setpoint (or posture) \mathbf{u} in the actuator space, that is supposed to be reached at short notice by the motor apparatus once the command is emitted, under classical stability/controllability constraints. This entails that, consistently with the ‘end-effector’ ballistic control setup (Mussa-Ivaldi and Solla, 2004), \mathbf{u} is independent from \mathbf{u}_0 , i.e.:

$$\mathbf{u} \sim \Pr(U | \mathbf{a})$$

The motor command \mathbf{a} then corresponds to the desired end-orientation of the sensor, here considered as a setpoint in the actuators space, either expressed in actuators or endpoint coordinates (with hardware-implemented detailed effector response function). Under that perspective, the effector acts on the sensors position and orientation so as to achieve a certain perspective (or view) over the external scene, and the controlled emitter \mathbf{u} is now called a *viewpoint*.

(iii) **Uncontrolled environment** The third important assumption is that the motor command \mathbf{a} is not expected to affect the uncontrolled latent emitter \mathbf{z} , i.e.

$$\mathbf{z} \sim \Pr(Z | \mathbf{z}_0)$$

so that \mathbf{z} should depend only on the external dynamics (the external “uncontrolled” process).

(iv) **Static assumption** Under a scene decoding task, it is rather common to consider the environment as “static” (Butko and Movellan, 2010). This fourth assumption means, in short, that:

$$\Pr(Z | \mathbf{z}_0) = \delta(Z, \mathbf{z}_0)$$

with δ the Kronecker symbol. The uncontrolled latent emitter \mathbf{z} is thus expected to capture all relevant information about the current scene, while remaining invariant throughout the decoding process.

Last, the observation \mathbf{x} may rely on both emitters \mathbf{z} and \mathbf{u} , i.e.

$$\mathbf{x} \sim \Pr(X | \mathbf{z}, \mathbf{u}) \tag{39}$$

Each observation \mathbf{x} is generated from a mixed emitter (\mathbf{z}, \mathbf{u}) , with \mathbf{u} the controlled part of the emitter and \mathbf{z} the uncontrolled part. Note that \mathbf{z} is said the latent state out of habit, though both \mathbf{u} and \mathbf{z} contribute to the generation of \mathbf{x} .

For notational simplicity, we absorb here the execution noise (Van Beers et al., 2004) in the measure process, i.e.: $\mathbf{x} \sim \Pr(X | \mathbf{z}, U) \Pr(U | \mathbf{a})$. Then, by notational abuse, we assimilate in the rest

of the paper \mathbf{u} (the controlled emitter) with \mathbf{a} (the motor command), so that a single variable $\mathbf{u} \equiv \mathbf{a}$ should be used for both. Each different \mathbf{u} is thus both interpreted as a motor command and as an emitter. As a motor command, it is controllable, i.e. determined by a controller. As an emitter, it monitors the generation of the sensory field, in combination with the latent state \mathbf{z} .

B Viewpoint-dependent variational encoding setup

The *variational encoding* perspective (Hinton and Zemel, 1994) was originally developed to train unsupervised autoencoder neural networks. If \mathbf{x} is the original data, the corresponding code \mathbf{z} is generated by a distribution q , i.e. $\mathbf{z} \sim q(Z)$. This distribution is called the *encoder*. Then, the reconstruction is made possible with a second conditional probability over the codes, i.e. $p(X|\mathbf{z})$, that is called the *decoder*. If \mathbf{z} is the current code, the reconstructed data is $\tilde{\mathbf{x}} \sim p(X|\mathbf{z})$.

In short, the efficacy of a code is estimated by an information-theoretic quantity, the “reconstruction cost” that is defined for every \mathbf{x} knowing p and q :

$$F(\mathbf{x}) = \mathbb{E}_{\mathbf{z} \sim q} [-\log(p(\mathbf{x}|\mathbf{z}))] + \text{KL}(q(Z)||p(Z)) \quad (40)$$

$$= -\log p(\mathbf{x}) + \text{KL}(q(Z)||p(Z|\mathbf{x})) \quad (41)$$

with $p(Z)$ the prior over the latent state. F is also said the Variational Free Energy (VFE), for it shares a mathematic analogy with the Helmholtz Free Energy (Friston, 2010). Minimizing the cost F according to p and q thus means minimizing the “surprise” caused by observing the data \mathbf{x} (Friston, 2010).

Viewpoint-dependent VFE If we now turn back to the viewpoint selection setup, an additional factor \mathbf{u} (the viewpoint) comes into the play. The data \mathbf{x} that is actually read is now conditioned on \mathbf{u} , so that:

$$F(\mathbf{x}|\mathbf{u}) = \mathbb{E}_{\mathbf{z} \sim q} [-\log(p(\mathbf{x}|\mathbf{z}, \mathbf{u}))] + \text{KL}(q(Z)||p(Z)) \quad (42)$$

$$= -\log p(\mathbf{x}|\mathbf{u}) + \text{KL}(q(Z)||p(Z|\mathbf{x}, \mathbf{u})) \quad (43)$$

When only the variations of p and q are considered in the optimization, each viewpoint \mathbf{u} provides a distinct optimization problem that is resolved by finding $q(Z) \simeq p(Z|\mathbf{x}, \mathbf{u})$. Each \mathbf{u} may thus drive a different posterior and thus a different reconstruction cost. It is thus feasible to change (and optimize) the reconstruction cost through changing \mathbf{u} .

Sequential viewpoint-dependent VFE When generalized to many observations: $(\mathbf{x}, \mathbf{u}), (\mathbf{x}', \mathbf{u}'), \dots, (\mathbf{x}^{(n)}, \mathbf{u}^{(n)})$, the n^{th} reconstruction cost $F^{(n)}(\mathbf{x}^{(n)}|\mathbf{u}^{(n)}, \dots, \mathbf{x}, \mathbf{u})$ also obeys to the chain rule (see eq. 6), i.e. is estimated from $q^{(n-1)}, \mathbf{u}^{(n)}$ and $\mathbf{x}^{(n)}$ only:

$$F(\mathbf{x}^{(n)}|\mathbf{u}^{(n)}; q^{(n-1)}) = \mathbb{E}_{\mathbf{z} \sim q} [-\log p(\mathbf{x}^{(n)}|\mathbf{z}, \mathbf{u}^{(n)})] + \text{KL}(q(Z)||q^{(n-1)}(Z)) \quad (44)$$

$$= -\log p(\mathbf{x}^{(n)}|\mathbf{u}^{(n)}) + \text{KL}(q(Z)||p(Z|\mathbf{x}^{(n)}, \mathbf{u}^{(n)}; q^{(n-1)})) \quad (45)$$

with $q^{(n-1)}$ having the role of the prior, providing a *forward* variational encoding scheme (see also (Chung et al., 2015; Fraccaro et al., 2016)).