

A Physical Account of the Free Energy Principle

Dalton A R Sakthivadivel

VERSES Research Lab

10th August 2023

A disclaimer

I am not Karl Friston

- ▶ I have my own ideas about what the FEP is — what it 'should' say and how it ought to say it
- ▶ In many cases my thoughts are equivalent to his up to mathematical nuance (but are my own)
- ▶ I will point out what is Friston's and what is mine, and what these similarities and differences are, explicitly

A silver lining

I can speak to many of Friston's ideas

And put them in the language of more traditional mathematics and physics

Today we will discuss one form of his theory which applies to non-equilibrium statistical mechanics

We will

- ▶ take some general ideas and concepts from Friston's work
- ▶ turn them into mathematical statements
- ▶ prove them (in sketches)
- ▶ and use these theorems in an explicit model from NESM

Synchronisation and approximate Bayesian inference

We begin from a system of two coupled random variables evolving in time separated by a boundary,

$$X_t \xrightarrow{g} B_t \xrightarrow{h} Y_t$$

assumed to satisfy Itô SDEs

$$dX_t = f_1(X_t, B_t, t) dt + \sqrt{D_1(X_t, B_t, t)} dW_t^1$$

$$dB_t = f_2(X_t, B_t, Y_t, t) dt + \sqrt{D_2(X_t, B_t, Y_t, t)} dW_t^2$$

$$dY_t = f_3(B_t, Y_t, t) dt + \sqrt{D_3(B_t, Y_t, t)} dW_t^3$$

The precise coupling structure is specific to a given system and defines different classes of dynamics [Friston et al, 2023].

Synchronisation and approximate Bayesian inference

We begin from a system of two coupled random variables evolving in time separated by a boundary,

$$X_t \xrightarrow{g} B_t \xrightarrow{h} Y_t$$

assumed to satisfy Itô SDEs

$$dX_t = f_1(X_t, B_t, t) dt + \sqrt{D_1(X_t, B_t, t)} dW_t^1$$

$$dB_t = f_2(X_t, B_t, Y_t, t) dt + \sqrt{D_2(X_t, B_t, Y_t, t)} dW_t^2$$

$$dY_t = f_3(B_t, Y_t, t) dt + \sqrt{D_3(B_t, Y_t, t)} dW_t^3$$

The precise coupling structure is specific to a given system and defines different classes of dynamics [Friston et al, 2023].

Synchronisation and approximate Bayesian inference

We begin from a system of two coupled random variables evolving in time separated by a boundary,

$$X_t \xrightarrow{g} B_t \xrightarrow{h} Y_t$$

assumed to satisfy Itô SDEs

$$dX_t = f_1(X_t, B_t, t) dt + \sqrt{D_1(X_t, B_t, t)} dW_t^1$$

$$dB_t = f_2(X_t, B_t, Y_t, t) dt + \sqrt{D_2(X_t, B_t, Y_t, t)} dW_t^2$$

$$dY_t = f_3(B_t, Y_t, t) dt + \sqrt{D_3(B_t, Y_t, t)} dW_t^3$$

The precise coupling structure is specific to a given system and defines different classes of dynamics [Friston et al, 2023].

Synchronisation continued

Now define a function u mapping conditional means to conditional means such that

$$\mathbf{E}[Y_t | B_t] = u(\mathbf{E}[X_t | B_t]).$$

More generally, suppose there exists a path-wise attractor for X_t dependent on B_t , (likewise for Y_t) such that u maps attracting states to attracting states

This utilises a version of centre manifold theory adapted to stochastic processes [Brzeźniak, Capiński, and Flandoli 1993]

Synchronisation continued

Now define a function u mapping conditional means to conditional means such that

$$\mathbf{E}[Y_t | B_t] = u(\mathbf{E}[X_t | B_t]).$$

More generally, suppose there exists a path-wise attractor for X_t dependent on B_t , (likewise for Y_t) such that u maps attracting states to attracting states

This utilises a version of centre manifold theory adapted to stochastic processes [Brzeźniak, Capiński, and Flandoli 1993]

Synchronisation continued

Now define a function u mapping conditional means to conditional means such that

$$\mathbf{E}[Y_t | B_t] = u(\mathbf{E}[X_t | B_t]).$$

More generally, suppose there exists a path-wise attractor for X_t dependent on B_t , (likewise for Y_t) such that u maps attracting states to attracting states

This utilises a version of centre manifold theory adapted to stochastic processes [Brzeźniak, Capiński, and Flandoli 1993]

Synchronisation continued

Y_t does not depend on X_t (only B_t)

Therefore X_t does not affect Y_t directly

But it can propagate changes to Y_t through B_t

By u , X_t parameterises a set of likely paths of Y_t : by construction, every Y_t must arise from the propagation of the influence of X_t through B_t

→ There is a likelihood of Y_t paths associated to every X_t path

Synchronisation continued

Y_t does not depend on X_t (only B_t)

Therefore X_t does not affect Y_t directly

But it can propagate changes to Y_t through B_t

By u , X_t parameterises a set of likely paths of Y_t : by construction, every Y_t must arise from the propagation of the influence of X_t through B_t

→ There is a likelihood of Y_t paths associated to every X_t path

Synchronisation continued

Y_t does not depend on X_t (only B_t)

Therefore X_t does not affect Y_t directly

But it can propagate changes to Y_t through B_t

By u , X_t parameterises a set of likely paths of Y_t : by construction, every Y_t must arise from the propagation of the influence of X_t through B_t

→ There is a likelihood of Y_t paths associated to every X_t path

Synchronisation continued

Y_t does not depend on X_t (only B_t)

Therefore X_t does not affect Y_t directly

But it can propagate changes to Y_t through B_t

By u , X_t parameterises a set of likely paths of Y_t : by construction, every Y_t must arise from the propagation of the influence of X_t through B_t

→ There is a likelihood of Y_t paths associated to every X_t path

Synchronisation continued

Y_t does not depend on X_t (only B_t)

Therefore X_t does not affect Y_t directly

But it can propagate changes to Y_t through B_t

By u , X_t parameterises a set of likely paths of Y_t : by construction, every Y_t must arise from the propagation of the influence of X_t through B_t

→ There is a likelihood of Y_t paths associated to every X_t path

Synchronisation continued

We now define the following two path probability densities:
 $q(y(t); m_{y,b}(t), \varrho_{y,b}(t))$ and $p(y(t) | b(t))$

q is a density of Y_t parameterised by conditional moments

p is the conditional density of Y_t derived from the joint distribution

We ask that these are the same

$$D_{KL}(q||p) \equiv 0 \quad (p\text{-a.s.})$$

Synchronisation continued

We now define the following two path probability densities:
 $q(y(t); m_{y,b}(t), \varrho_{y,b}(t))$ and $p(y(t) | b(t))$

q is a density of Y_t parameterised by conditional moments

p is the conditional density of Y_t derived from the joint distribution

We ask that these are the same

$$D_{KL}(q||p) \equiv 0 \quad (p\text{-a.s.})$$

Synchronisation continued

We now define the following two path probability densities:
 $q(y(t); m_{y,b}(t), \varrho_{y,b}(t))$ and $p(y(t) | b(t))$

q is a density of Y_t parameterised by conditional moments

p is the conditional density of Y_t derived from the joint distribution

We ask that these are the same

$$D_{KL}(q||p) \equiv 0 \quad (p\text{-a.s.})$$

Synchronisation continued

We now define the following two path probability densities:
 $q(y(t); m_{y,b}(t), \varrho_{y,b}(t))$ and $p(y(t) | b(t))$

q is a density of Y_t parameterised by conditional moments

p is the conditional density of Y_t derived from the joint distribution

We ask that these are the same

$$D_{KL}(q||p) \equiv 0 \quad (p\text{-a.s.})$$

Synchronisation continued

We now define the following two path probability densities:
 $q(y(t); m_{y,b}(t), \varrho_{y,b}(t))$ and $p(y(t) | b(t))$

q is a density of Y_t parameterised by conditional moments

p is the conditional density of Y_t derived from the joint distribution

We ask that these are the same

$$D_{KL}(q||p) \equiv 0 \quad (p\text{-a.s.})$$

Synchronisation continued

By our previous argument, if the coupling exists and maps conditional means to conditional means, we can rewrite q as

$$q(y(t); u(x(t)), \varrho_{y,b}(t))$$

for some choice of $X_t = x(t)$

Each choice will produce a different likelihood of Y_t paths; only one choice of $x(t)$ produces the density that matches q to p

By direct substitution we verify that $q = p$ (almost surely) when $x(t)$ equals the expected $x(t)$ given $b(t)$ and the variance of $x(t)$ is the same as the variance of $y(t)$ given $b(t)$

So if we want $D_{KL} = 0$, then we must see “conditional mode-matching”

Synchronisation continued

By our previous argument, if the coupling exists and maps conditional means to conditional means, we can rewrite q as

$$q(y(t); u(x(t)), \varrho_{y,b}(t))$$

for some choice of $X_t = x(t)$

Each choice will produce a different likelihood of Y_t paths; only one choice of $x(t)$ produces the density that matches q to p

By direct substitution we verify that $q = p$ (almost surely) when $x(t)$ equals the expected $x(t)$ given $b(t)$ and the variance of $x(t)$ is the same as the variance of $y(t)$ given $b(t)$

So if we want $D_{KL} = 0$, then we must see “conditional mode-matching”

Synchronisation continued

By our previous argument, if the coupling exists and maps conditional means to conditional means, we can rewrite q as

$$q(y(t); u(x(t)), \varrho_{y,b}(t))$$

for some choice of $X_t = x(t)$

Each choice will produce a different likelihood of Y_t paths; only one choice of $x(t)$ produces the density that matches q to p

By direct substitution we verify that $q = p$ (almost surely) when $x(t)$ equals the expected $x(t)$ given $b(t)$ and the variance of $x(t)$ is the same as the variance of $y(t)$ given $b(t)$

So if we want $D_{KL} = 0$, then we must see “conditional mode-matching”

Synchronisation continued

By our previous argument, if the coupling exists and maps conditional means to conditional means, we can rewrite q as

$$q(y(t); u(x(t)), \varrho_{y,b}(t))$$

for some choice of $X_t = x(t)$

Each choice will produce a different likelihood of Y_t paths; only one choice of $x(t)$ produces the density that matches q to p

By direct substitution we verify that $q = p$ (almost surely) when $x(t)$ equals the expected $x(t)$ given $b(t)$ and the variance of $x(t)$ is the same as the variance of $y(t)$ given $b(t)$

So if we want $D_{KL} = 0$, then we must see “conditional mode-matching”

Synchronisation and approximate Bayesian inference

So far all we have said is that, *via* the interactions across a shared boundary, coupled random dynamical systems estimate each others statistics in a very literal sense

Example: any 'thing' encodes a probability distribution over possible environmental states because the environment must be conducive to it existing [S 2022; Ramstead, S et al 2023]

Synchronisation and approximate Bayesian inference

So far all we have said is that, *via* the interactions across a shared boundary, coupled random dynamical systems estimate each others statistics in a very literal sense

Example: any 'thing' encodes a probability distribution over possible environmental states because the environment must be conducive to it existing [S 2022; Ramstead, S et al 2023]

Synchronisation continued

Bayesian inference is parameter estimation, meaning we have statements about approximate Bayesian inference (only inferring two parameters and only making inferences about the environment instead of the object-boundary-environment system)

Complex systems are hard to understand because of their interactions, so replacing this with the study of variational free energy is fruitful. Experimental support for this claim is available [Isomura et al 2023]

Synchronisation continued

Bayesian inference is parameter estimation, meaning we have statements about approximate Bayesian inference (only inferring two parameters and only making inferences about the environment instead of the object-boundary-environment system)

Complex systems are hard to understand because of their interactions, so replacing this with the study of variational free energy is fruitful. Experimental support for this claim is available [Isomura et al 2023]

An inequality relating to organisation

Since $VFE \geq -\log p(x(t), b(t))$, minimising free energy optimises an upper bound on surprisal

In effect we are saying that if the system mainly does what we expect it to do, it can only be so surprising (for instance, stones must be concentrated on stone-like states; control systems must be concentrated on set points)

An explicit model

Consider a ferromagnet m_t , a heat bath B_t , and a source or sink S_t (all in the presence of noise)

S_t acts on m_t through B_t (supplying or subtracting heat to the bath in contact with m_t)

Let $S > 0$ be a source state (<, sink, resp)

Suppose the initial temperature of B was $T > T_c$. If the expected $m(t)$ given $B(t)$ reaches ± 1 at some t then it is likely some heat was pumped out (so $S(t) < 0$ on some interval). The inverse is also true

So $q(S(t); u(m_t)) \sim p(S(t) | B(t))$ when $m(t)$ is the mean state paired with the specific value of $B(t)$.

An explicit model

Consider a ferromagnet m_t , a heat bath B_t , and a source or sink S_t (all in the presence of noise)

S_t acts on m_t through B_t (supplying or subtracting heat to the bath in contact with m_t)

Let $S > 0$ be a source state (<, sink, resp)

Suppose the initial temperature of B was $T > T_c$. If the expected $m(t)$ given $B(t)$ reaches ± 1 at some t then it is likely some heat was pumped out (so $S(t) < 0$ on some interval). The inverse is also true

So $q(S(t); u(m_t)) \sim p(S(t) | B(t))$ when $m(t)$ is the mean state paired with the specific value of $B(t)$.

An explicit model

Consider a ferromagnet m_t , a heat bath B_t , and a source or sink S_t (all in the presence of noise)

S_t acts on m_t through B_t (supplying or subtracting heat to the bath in contact with m_t)

Let $S > 0$ be a source state (<, sink, resp)

Suppose the initial temperature of B was $T > T_c$. If the expected $m(t)$ given $B(t)$ reaches ± 1 at some t then it is likely some heat was pumped out (so $S(t) < 0$ on some interval). The inverse is also true

So $q(S(t); u(m_t)) \sim p(S(t) | B(t))$ when $m(t)$ is the mean state paired with the specific value of $B(t)$.

An explicit model

Consider a ferromagnet m_t , a heat bath B_t , and a source or sink S_t (all in the presence of noise)

S_t acts on m_t through B_t (supplying or subtracting heat to the bath in contact with m_t)

Let $S > 0$ be a source state (<, sink, resp)

Suppose the initial temperature of B was $T > T_c$. If the expected $m(t)$ given $B(t)$ reaches ± 1 at some t then it is likely some heat was pumped out (so $S(t) < 0$ on some interval). The inverse is also true

So $q(S(t); u(m_t)) \sim p(S(t) | B(t))$ when $m(t)$ is the mean state paired with the specific value of $B(t)$.

An explicit model

Consider a ferromagnet m_t , a heat bath B_t , and a source or sink S_t (all in the presence of noise)

S_t acts on m_t through B_t (supplying or subtracting heat to the bath in contact with m_t)

Let $S > 0$ be a source state (<, sink, resp)

Suppose the initial temperature of B was $T > T_c$. If the expected $m(t)$ given $B(t)$ reaches ± 1 at some t then it is likely some heat was pumped out (so $S(t) < 0$ on some interval). The inverse is also true

So $q(S(t); u(m_t)) \sim p(S(t) | B(t))$ when $m(t)$ is the mean state paired with the specific value of $B(t)$.

Another explicit model

Give m_t access to a heat pump and assume this parameterisation result holds

The system can now make changes in the environment such that its preferred state is the 'right' parameter

This is sometimes called *active inference*

Slogans like "systems enact their own preferences" are sometimes invoked

Another explicit model

Give m_t access to a heat pump and assume this parameterisation result holds

The system can now make changes in the environment such that its preferred state is the 'right' parameter

This is sometimes called *active inference*

Slogans like "systems enact their own preferences" are sometimes invoked

Another explicit model

Give m_t access to a heat pump and assume this parameterisation result holds

The system can now make changes in the environment such that its preferred state is the 'right' parameter

This is sometimes called *active inference*

Slogans like "systems enact their own preferences" are sometimes invoked

Another explicit model

Give m_t access to a heat pump and assume this parameterisation result holds

The system can now make changes in the environment such that its preferred state is the 'right' parameter

This is sometimes called *active inference*

Slogans like “systems enact their own preferences” are sometimes invoked

Why thermodynamics

Despite the name, FEP is *not* primarily a theory of thermodynamics (even in Friston's work it is closer to control theory than statistical physics)

Despite the prevailing conversations, I think of the mathematics of the FEP as about cohesive things in heat baths — *i.e.* not necessarily self-organisation (and not necessarily brains). It is strictly more general

The point is to reduce our characterisation of coupled random dynamical systems to Bayesian statistics

It does however fit into the thermodynamics of self-organisation

Why thermodynamics

Despite the name, FEP is *not* primarily a theory of thermodynamics (even in Friston's work it is closer to control theory than statistical physics)

Despite the prevailing conversations, I think of the mathematics of the FEP as about cohesive things in heat baths — *i.e.* not necessarily self-organisation (and not necessarily brains). It is strictly more general

The point is to reduce our characterisation of coupled random dynamical systems to Bayesian statistics

It does however fit into the thermodynamics of self-organisation

Why thermodynamics

Despite the name, FEP is *not* primarily a theory of thermodynamics (even in Friston's work it is closer to control theory than statistical physics)

Despite the prevailing conversations, I think of the mathematics of the FEP as about cohesive things in heat baths — *i.e.* not necessarily self-organisation (and not necessarily brains). It is strictly more general

The point is to reduce our characterisation of coupled random dynamical systems to Bayesian statistics

It does however fit into the thermodynamics of self-organisation

Why thermodynamics

Despite the name, FEP is *not* primarily a theory of thermodynamics (even in Friston's work it is closer to control theory than statistical physics)

Despite the prevailing conversations, I think of the mathematics of the FEP as about cohesive things in heat baths — *i.e.* not necessarily self-organisation (and not necessarily brains). It is strictly more general

The point is to reduce our characterisation of coupled random dynamical systems to Bayesian statistics

It does however fit into the thermodynamics of self-organisation

Thermodynamics

Suppose a system minimises free energy. What can we say about its interactions with the environment?

Because there are interactions we can say something about flows in and out of the system

If the flows have certain properties there exists a NESS (physics statement), and if the system minimises variational FE there exists a NESS (control theory statement)

What's the relation?

Thermodynamics

Suppose a system minimises free energy. What can we say about its interactions with the environment?

Because there are interactions we can say something about flows in and out of the system

If the flows have certain properties there exists a NESS (physics statement), and if the system minimises variational FE there exists a NESS (control theory statement)

What's the relation?

Thermodynamics

Suppose a system minimises free energy. What can we say about its interactions with the environment?

Because there are interactions we can say something about flows in and out of the system

If the flows have certain properties there exists a NESS (physics statement), and if the system minimises variational FE there exists a NESS (control theory statement)

What's the relation?

Thermodynamics

Suppose a system minimises free energy. What can we say about its interactions with the environment?

Because there are interactions we can say something about flows in and out of the system

If the flows have certain properties there exists a NESS (physics statement), and if the system minimises variational FE there exists a NESS (control theory statement)

What's the relation?

Thermodynamics

Suppose a system minimises free energy. What can we say about its interactions with the environment?

Because there are interactions we can say something about flows in and out of the system

If the flows have certain properties there exists a NESS (physics statement), and if the system minimises variational FE there may exist a NESS (more like a control theory statement)

What's the relation?

Thermodynamics

Suppose a system minimises free energy. What can we say about its interactions with the environment?

Because there are interactions we can say something about flows in and out of the system

If the flows have certain properties there exists a NESS (physics statement), and if the system minimises variational FE there may exist a NESS (more like a control theory statement)

What's the relation?

Thermodynamics

Suppose a system minimises free energy. What can we say about its interactions with the environment?

Because there are interactions we can say something about flows in and out of the system

If the flows have certain properties there exists a NESS (physics statement), and if the system minimises variational FE there may exist a NESS (more like a control theory statement)

What's the relation?

Thermodynamics

Suppose a system minimises free energy. What can we say about its interactions with the environment?

Because there are interactions we can say something about flows in and out of the system

If the flows have certain properties there exists a NESS (physics statement), and if the system minimises variational FE there may exist a NESS (more like a control theory statement)

What's the relation?

Thermodynamics of self-organisation in active inference

In our example we saw how the FEP entails an object forming a model of the fluxes of heat through it and its environment

Due to environmental perturbations, there is an energetic cost to maintaining a NESS

There are energetic resources in the environment

A good model of an environment allows the system to (i) find resources (ii) track perturbations (iii) continue to exist

A sketch of this argument first appears in [Ueltzhöffer 2020]

Thermodynamics of self-organisation in active inference

In our example we saw how the FEP entails an object forming a model of the fluxes of heat through it and its environment

Due to environmental perturbations, there is an energetic cost to maintaining a NESS

There are energetic resources in the environment

A good model of an environment allows the system to (i) find resources (ii) track perturbations (iii) continue to exist

A sketch of this argument first appears in [Uetzhöffer 2020]

Thermodynamics of self-organisation in active inference

In our example we saw how the FEP entails an object forming a model of the fluxes of heat through it and its environment

Due to environmental perturbations, there is an energetic cost to maintaining a NESS

There are energetic resources in the environment

A good model of an environment allows the system to (i) find resources (ii) track perturbations (iii) continue to exist

A sketch of this argument first appears in [Ueltzhöffer 2020]

Thermodynamics of self-organisation in active inference

In our example we saw how the FEP entails an object forming a model of the fluxes of heat through it and its environment

Due to environmental perturbations, there is an energetic cost to maintaining a NESS

There are energetic resources in the environment

A good model of an environment allows the system to (i) find resources (ii) track perturbations (iii) continue to exist

A sketch of this argument first appears in [Ueltzhöffer 2020]

Thermodynamics of self-organisation in active inference

In our example we saw how the FEP entails an object forming a model of the fluxes of heat through it and its environment

Due to environmental perturbations, there is an energetic cost to maintaining a NESS

There are energetic resources in the environment

A good model of an environment allows the system to (i) find resources (ii) track perturbations (iii) continue to exist

A sketch of this argument first appears in [Ueltzhöffer 2020]

Thermodynamics of self-organisation in active inference

The takeaway: systems with good models of the physics of their environments are better at self-organising; conversely systems which are good at self-organising are good sources of information about what the environments they are in must look like (this direction is more in the spirit of the FEP)

Thermodynamics of self-organisation in active inference

The takeaway: systems with good models of the physics of their environments are better at self-organising; conversely systems which are good at self-organising are good sources of information about what the environments they are in must look like (this direction is more in the spirit of the FEP)

Max cal?

Due to its connections to quantities on paths and path probability densities it is hypothesised that the FEP has a nice story to tell in the context of max cal, which can be used in interacting or non-equilibrium settings [Jaynes 1980]

Some connections to thermodynamics have already been written [Parr, Da Costa, and Friston 2020] as well as some connections to path integral formalisms and physics [S 2022]

Max cal?

Due to its connections to quantities on paths and path probability densities it is hypothesised that the FEP has a nice story to tell in the context of max cal, which can be used in interacting or non-equilibrium settings [Jaynes 1980]

Some connections to thermodynamics have already been written [Parr, Da Costa, and Friston 2020] as well as some connections to path integral formalisms and physics [S 2022]

- [1] K J Friston, L Da Costa, D A R Sakthivadivel, C Heins, G A Pavliotis, M Ramstead, and T Parr. *Phys Life Rev* (2023). Forthcoming.
- [2] Z Brzeźniak, M Capiński, and F Flandoli. *Probab Theory Relat Fields* **95**, 1 (1993).
- [3] D A R Sakthivadivel, *Active Inference: Third International Workshop* (Grenoble, 2022).
- [4] M Ramstead, D A R Sakthivadivel, C Heins, M Koudahl, B Millidge, L Da Costa, B Klein, K J Friston. *Interface Focus* **13**, 3 (2023).
- [5] T Isomura, K Kotani, Y Jimbo, and K J Friston. *Nat Commun* **14** (2023).
- [6] K Ueltzhöffer. arXiv preprint (2020).
- [7] E T Jaynes. *Annu Rev Phys Chem* **31**, 1 (1980).
- [8] T Parr, L Da Costa, and K J Friston. *Philos Trans R Soc A* **378**, 2164 (2020).