

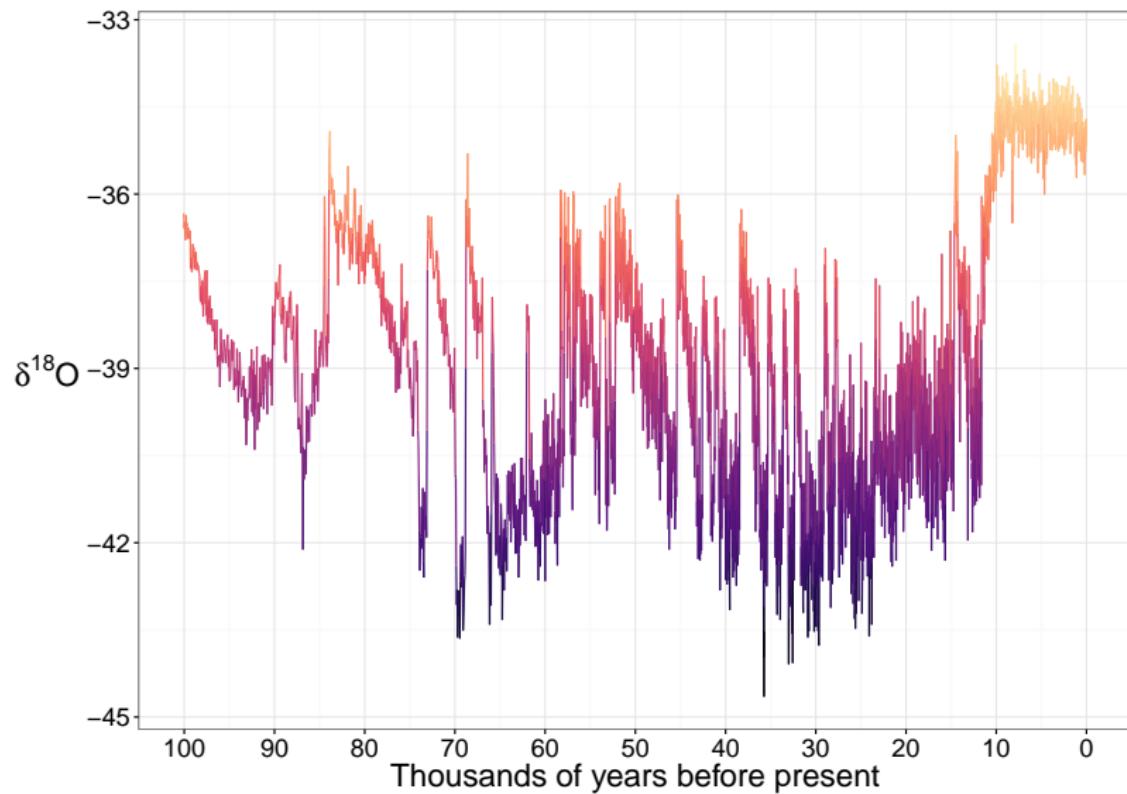
Statistical palaeoclimate reconstruction: how fast can climate change?

Andrew Parnell

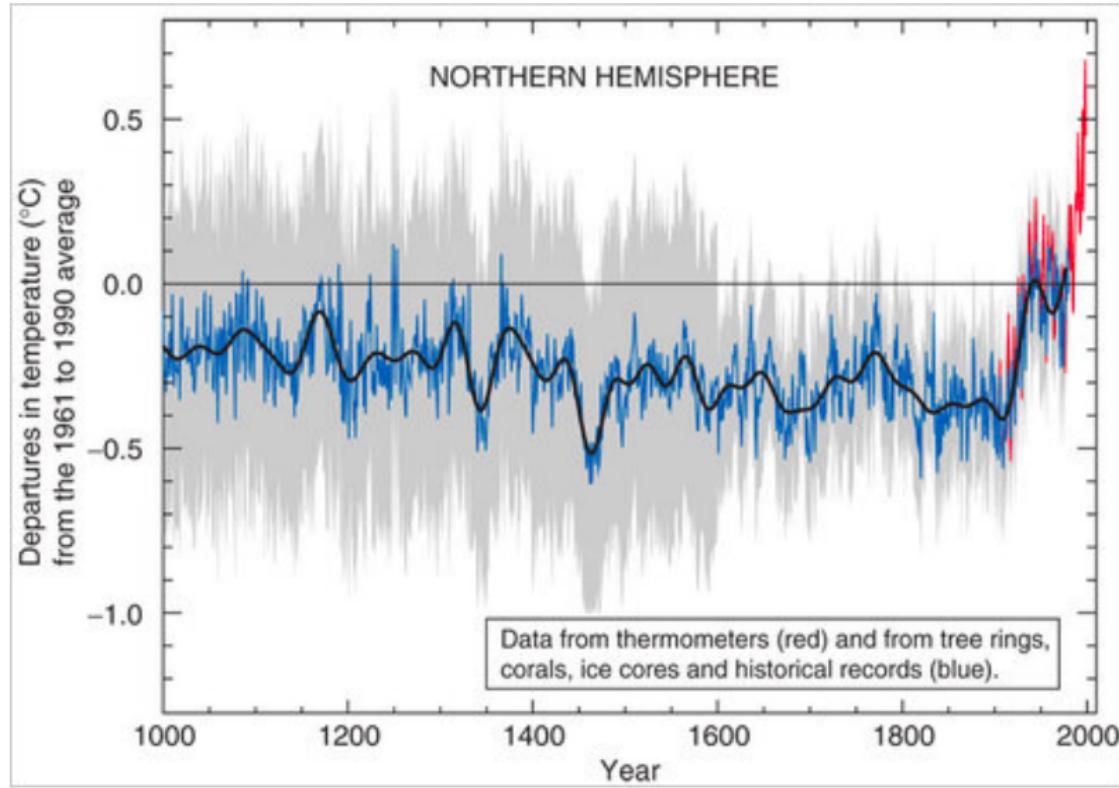
andrew.parnell@ucd.ie



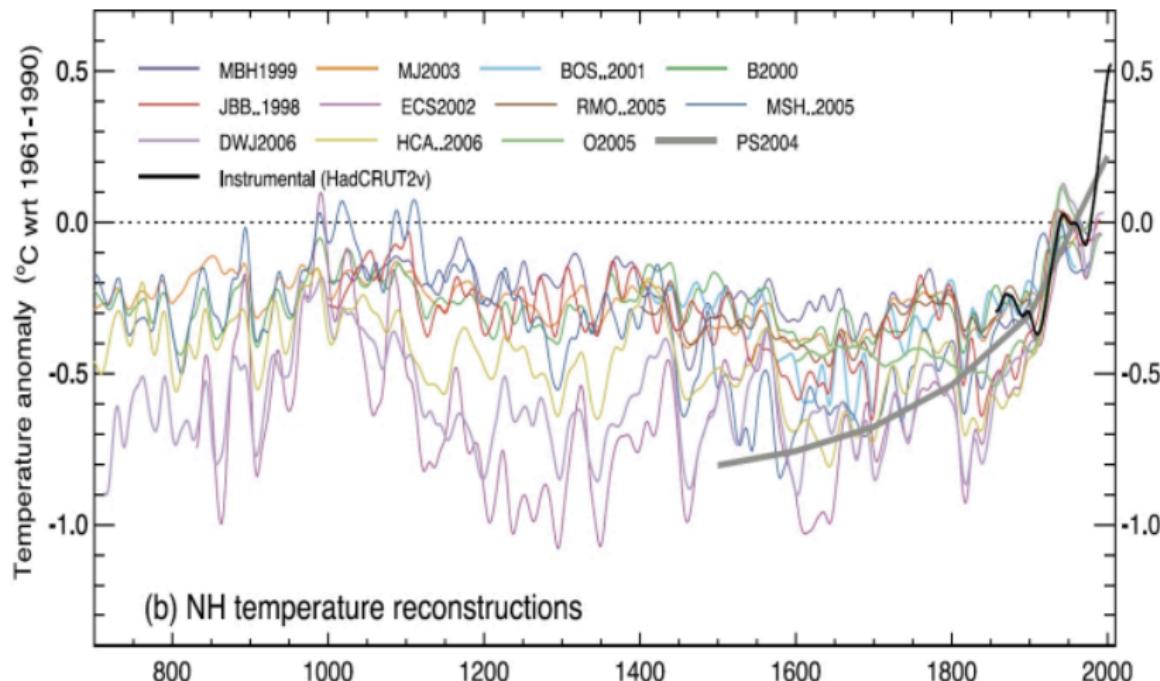
The GISP 2 ice core



The Hockey Stick



A bundle of Hockey Sticks



A Lago Grando di Monticchio



(By Pitichinaccio - Own work, Public Domain)

Cape May, New Jersey

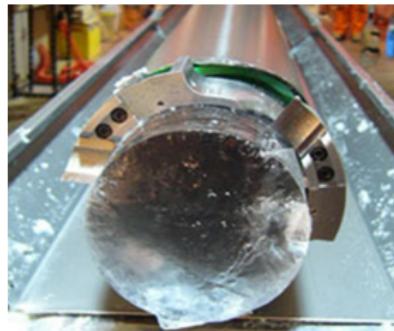


(By Smallbones - Own work, CC0)

Different proxies



(a) Tree rings



(b) An ice core

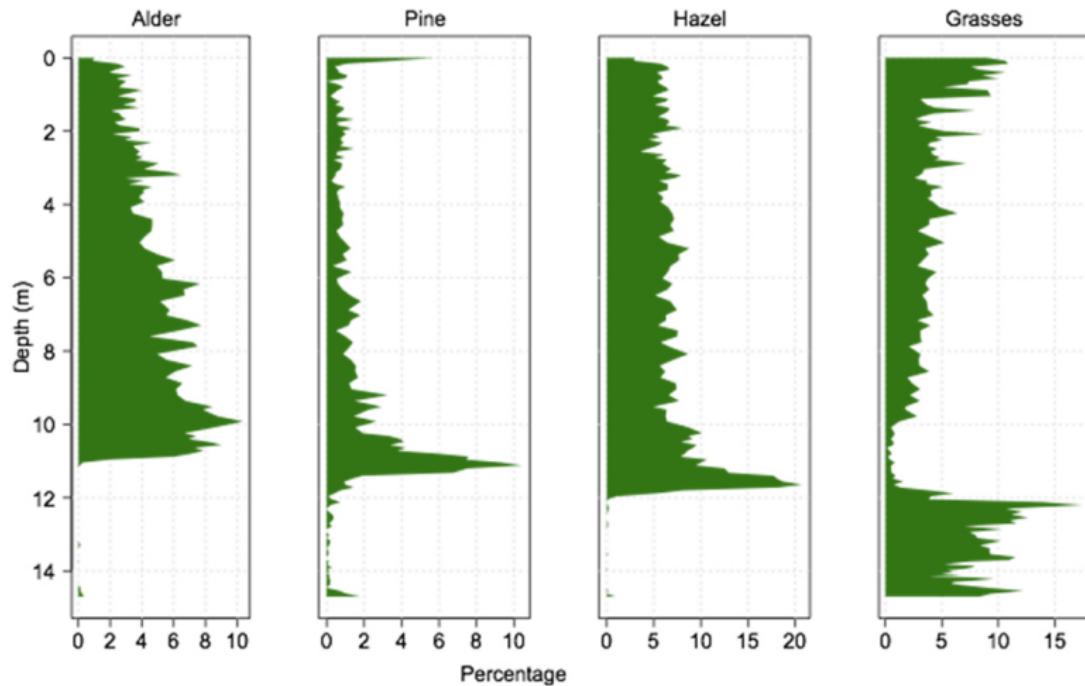


(c) A non-biting midge (chironomid)



(d) Alder pollen

Pollen depth plots



From pictures to data

Year	Climate	Proxy data			
2016	climate ₂₀₁₆	proxy _{1,2016}	proxy _{2,2016}	...	proxy _{p,2016}
2015	climate ₂₀₁₅	proxy _{1,2015}	proxy _{2,2015}	...	proxy _{p,2015}
:	:	:	:	...	:
1850	climate ₁₈₅₀	proxy _{1,1850}	proxy _{2,1850}	...	proxy _{p,1850}
1849	climate ₁₈₄₉	proxy _{1,1849}	proxy _{2,1849}	...	proxy _{p,1849}
:	:	:	:	...	:
1001	climate ₁₀₀₁	proxy _{1,1001}	proxy _{2,1001}	...	proxy _{p,1001}
1000	climate ₁₀₀₀	proxy _{1,1000}	proxy _{2,1000}	...	proxy _{p,1000}

A more general version

Calibration data set:

ID	Climate	Proxy data			
1	climate ₁	proxy _{1,1}	proxy _{2,1}	...	proxy _{p,1}
2	climate ₂	proxy _{1,2}	proxy _{2,2}	...	proxy _{p,2}
:	:	:	:	...	:
k	climate _k	proxy _{1,k}	proxy _{2,k}	...	proxy _{p,k}

Fossil data set:

Year	Climate	Proxy data			
n-1	climate _{n-1}	proxy _{1,n-1}	proxy _{2,n-1}	...	proxy _{p,n-1}
:	:	:	:	...	:
m+1	climate _{m+1}	proxy _{m+1}	proxy _{m+1}	...	proxy _{p,m+1}
m	climate _m	proxy _{1,m}	proxy _{2,m}	...	proxy _{p,m}

Some notation

Let:

- ▶ y be the ancient proxy data. Time indexed and usually multivariate
- ▶ c be ancient 'climate'. Time indexed and occasionally multivariate. Sometimes spatial too
- ▶ y^{cal} be the proxy data for the calibration period
- ▶ c^{cal} be the climate data for the calibration period

Aim is to find $c|y, y^{\text{cal}}, c^{\text{cal}}$

The regression version

Write:

$$c^{\text{cal}} = f(y^{\text{cal}}) + \epsilon$$

f might be a linear regression or involve some dimension reduction or variable selection.

The regression version

Write:

$$c^{\text{cal}} = f(y^{\text{cal}}) + \epsilon$$

f might be a linear regression or involve some dimension reduction or variable selection.

Then create:

$$\hat{c} = \hat{f}(y)$$

Problem solved!

Problems with this approach

Statistical:

- ▶ Hard to do model checking on f due to the size and nature of the calibration data
- ▶ c is often multivariate so people often pick one dimension
- ▶ The calibration period may be autocorrelated, leading to many spurious relationships
- ▶ Dimension reduction approaches will be very sensitive to the number of components chosen
- ▶ Lots of missing proxy data

Problems with this approach

Statistical:

- ▶ Hard to do model checking on f due to the size and nature of the calibration data
- ▶ c is often multivariate so people often pick one dimension
- ▶ The calibration period may be autocorrelated, leading to many spurious relationships
- ▶ Dimension reduction approaches will be very sensitive to the number of components chosen
- ▶ Lots of missing proxy data

Biological:

- ▶ The causation is the wrong way round. **Changes in climate cause changes in proxy values**
- ▶ The uncertainty in the proxies is usually substantial and not included
- ▶ The proxies might not be sensitive to northern hemisphere annual temperature, or any other chosen aspect of climate

A better way?

Instead write:

$$y^{\text{cal}} = f(c^{\text{cal}}) + \epsilon$$

f is known here as a **forward model** since it works in the causal direction. We can now include physical knowledge of how climate affects the proxies

A better way?

Instead write:

$$y^{\text{cal}} = f(c^{\text{cal}}) + \epsilon$$

f is known here as a **forward model** since it works in the causal direction. We can now include physical knowledge of how climate affects the proxies

Now **use Bayes**:

$$p(c|y, y^{\text{cal}}, c^{\text{cal}}) \propto p(y^{\text{cal}}|c^{\text{cal}})p(y|c)p(c)$$

We have the extra advantage that we can include a prior distribution $p(c)$ on the climate process

Bayesian palaeoclimate reconstruction in more detail

$$p(c, \theta, \phi | y, y^{\text{cal}}, c^{\text{cal}}) \propto p(y^{\text{cal}} | c^{\text{cal}}, \theta) p(y | c, \theta) p(c | \phi) p(\theta, \phi)$$

- ▶ $p(\theta, \phi)$ is a prior on the parameters that control the proxy/climate relationship, and climate dynamics respectively
- ▶ $p(c|\phi)$ is a prior distribution on climate dynamics. This might be a simple statistical time series model (e.g. a random walk) all the way up to a full general circulation model
- ▶ $p(y|c, \theta)$ is the forward model again, but this time applied to the missing ancient climates
- ▶ $p(y^{\text{cal}} | c^{\text{cal}}, \theta)$ is the forward model applied to the calibration data.

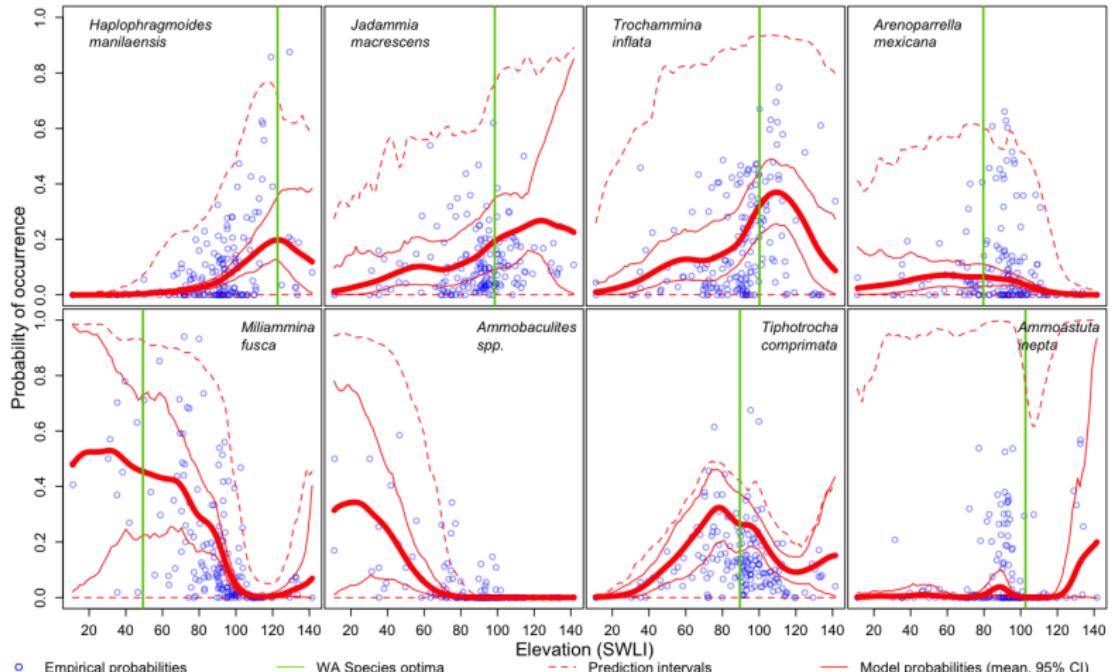
Why is this not the standard way people do this?

1. Building forward models is hard because you need a good calibration data set, some statistical modelling knowledge (especially with multivariate data), and some knowledge of the pollen/climate relationship
2. People want to avoid testing their models (out of sample evaluation etc)
3. Finding a good prior for climate dynamics is hard, especially if you have timing uncertainty
4. Bayes is still not common in climate science
5. Fitting this model to large calibration data sets is hard

Example: sea level rise in East Coast USA

- ▶ **Foramnifera** (or forams) live in the tidal range along coastal marshes
- ▶ There are lots of different species, and they all like slightly different parts of the tidal range
- ▶ If you take a sediment core on the marsh you can count lots of fossilised forams (which can also be dated) and produce a history of sea level height at that site
- ▶ We also take a number of surface samples from the local region to build up a calibration data set of the forams' preference for different aspects of the tidal range

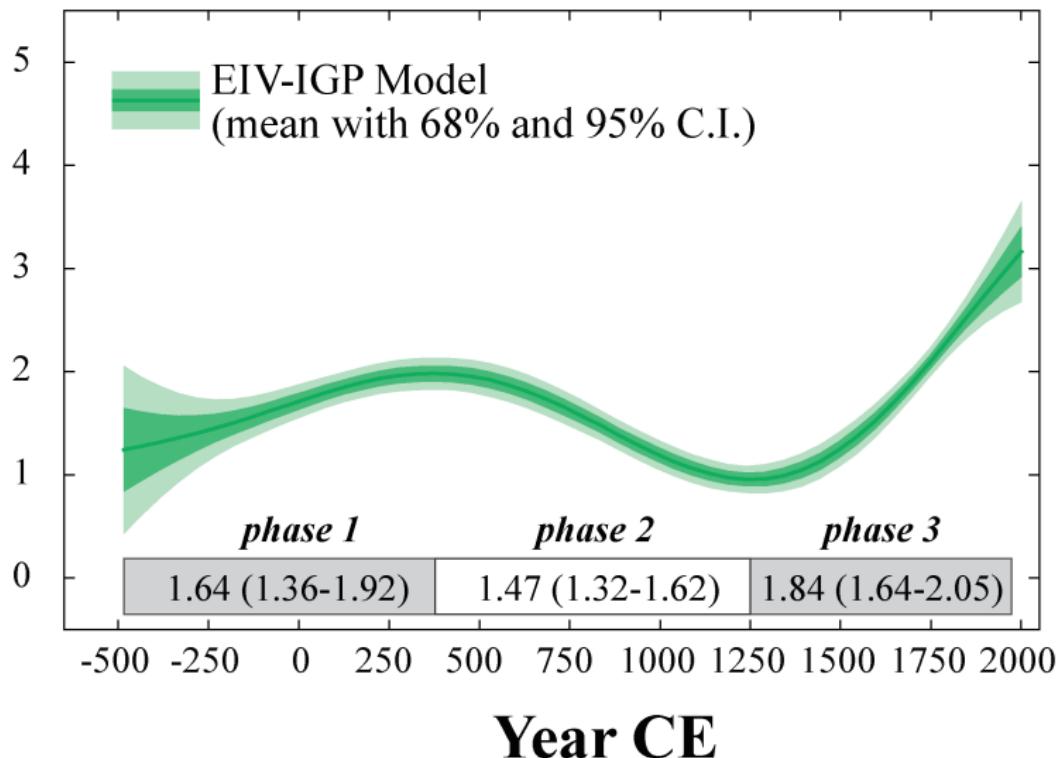
The forward model



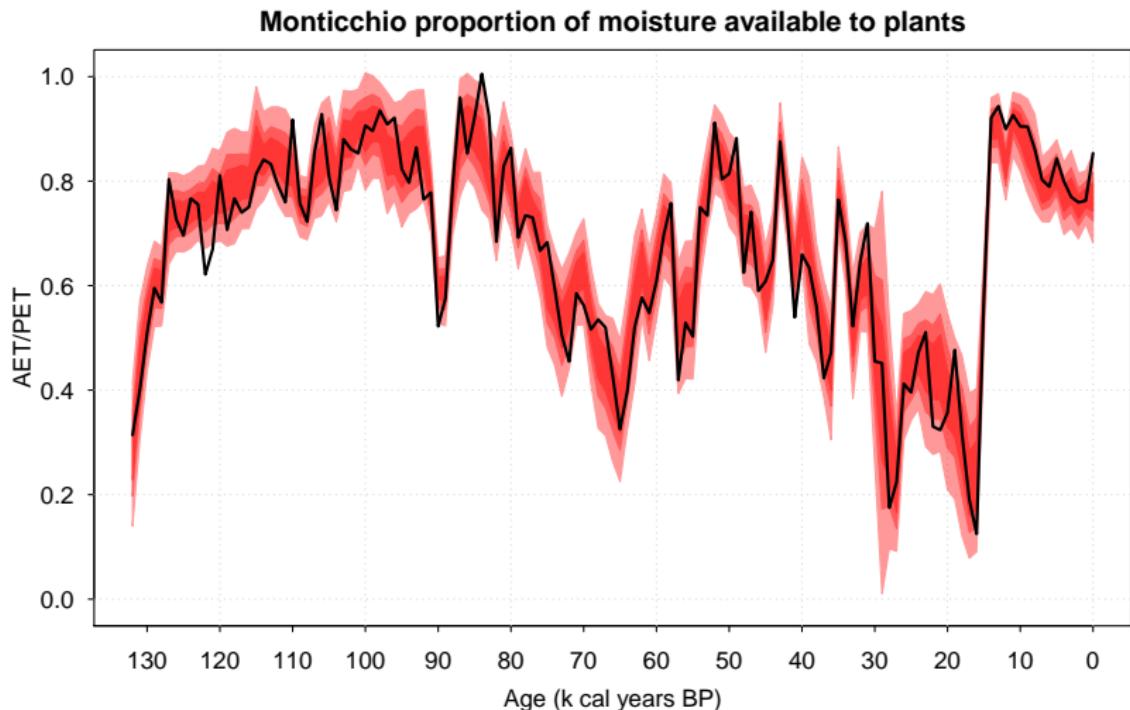
Model description

- ▶ Our forward model for the forams uses multinomial counts and P-splines
- ▶ We have a second proxy (called $\delta^{13}C$) that gives further information on the position in the tidal frame at that depth in the core
- ▶ Our prior on climate dynamics (here height of sea level over time) uses a fancy Gaussian process

Rate of sea level rise (mm/yr) for New Jersey, USA



Example 2: multivariate climate in Italy



Pollen forward models

- The forward model here is much more complicated, as we have ~10,000 modern pollen samples, with 3 dimensional climate and 28 compositional counts of pollen:

$$[y_1, \dots, y_{28}] \sim Mult(N, \{p_1, \dots, p_{28}\})$$

where, e.g.

$$p_I = \frac{\exp(\theta_I(c))}{\sum_j \exp(\theta_j(c))}$$

Pollen forward models

- ▶ The forward model here is much more complicated, as we have ~10,000 modern pollen samples, with 3 dimensional climate and 28 compositional counts of pollen:

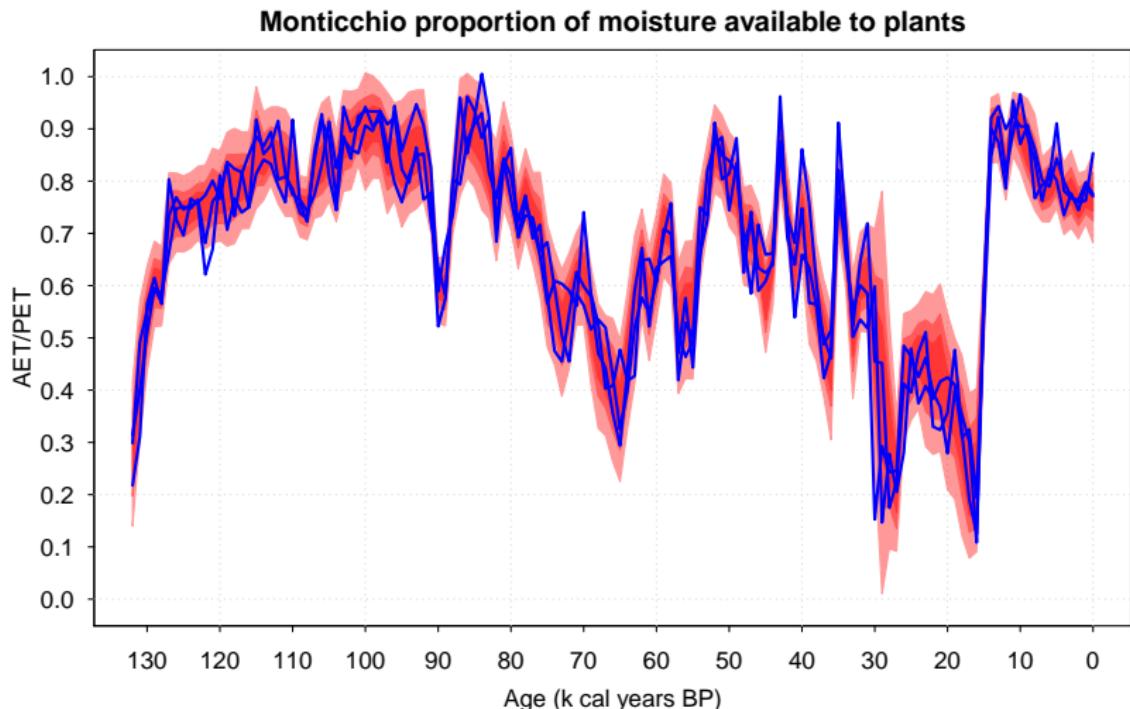
$$[y_1, \dots, y_{28}] \sim Mult(N, \{p_1, \dots, p_{28}\})$$

where, e.g.

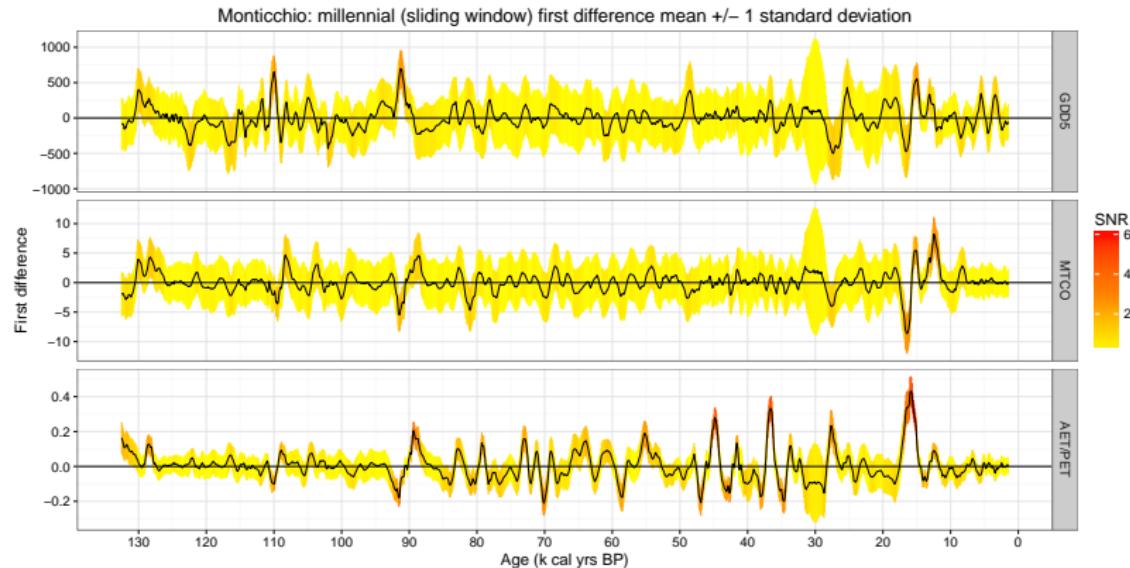
$$p_I = \frac{\exp(\theta_I(c))}{\sum_j \exp(\theta_j(c))}$$

- To fit such a model we make a small approximation and break the model into two, and fit the calibration parameters θ separately from the climate process parameters ϕ
- ▶ This places some mild restrictions on the climate process to still obtain valid inference

Example 2: Histories



Example 2: First differences - the speed of climate change



The grand challenge!

Fit a Bayesian model to:

- Reconstruct spatio-temporal palaeoclimate ...

The grand challenge!

Fit a Bayesian model to:

- Reconstruct spatio-temporal palaeoclimate ...
- ... using physical/statistical forward models for many proxies

The grand challenge!

Fit a Bayesian model to:

- Reconstruct spatio-temporal palaeoclimate ...
- ... using physical/statistical forward models for many proxies
- ... and physical/statistical models for climate dynamics

The grand challenge!

Fit a Bayesian model to:

- Reconstruct spatio-temporal palaeoclimate ...
- ... using physical/statistical forward models for many proxies
- ... and physical/statistical models for climate dynamics

The resulting output should be a large sample of spatio-temporal climate histories

Challenges 1: fitting state space models to large and complex data sets

What we really have is an **externally calibrated** state-space model in continuous time:

$$\begin{aligned}y^{\text{cal}}(t) &= f_{\theta}(c^{\text{cal}}(t)) + \epsilon^{\text{cal}} \\y(t) &= f_{\theta}(c(t)) + \epsilon\end{aligned}$$

Challenges 1: fitting state space models to large and complex data sets

What we really have is an **externally calibrated** state-space model in continuous time:

$$y^{\text{cal}}(t) = f_{\theta}(c^{\text{cal}}(t)) + \epsilon^{\text{cal}}$$

$$y(t) = f_{\theta}(c(t)) + \epsilon$$

$$c(t) = c_{\phi}(t - \Delta) + \gamma$$

Challenges 1: fitting state space models to large and complex data sets

What we really have is an **externally calibrated** state-space model in continuous time:

$$\begin{aligned}y^{\text{cal}}(t) &= f_{\theta}(c^{\text{cal}}(t)) + \epsilon^{\text{cal}} \\y(t) &= f_{\theta}(c(t)) + \epsilon \\c(t) &= c_{\phi}(t - \Delta) + \gamma\end{aligned}$$

- ▶ Fitting these models is hard when all the quantities are multivariate and f is a complex function
- ▶ Pseudo-marginal particle approaches seem to be the way to go for single-site models
- ▶ No obvious method yet for multi-site models. Perhaps SPDE-INLA?

Challenges 2: Incorporating mechanistic models

A more complex version (ignoring external calibration):

$$\begin{aligned}y(t) &= f_{\theta}(c(t)) \\c(t) &= g_{\phi}(c(t_-))\end{aligned}$$

Challenges 2: Incorporating mechanistic models

A more complex version (ignoring external calibration):

$$\begin{aligned}y(t) &= f_\theta(c(t)) \\c(t) &= g_\phi(c(t_-))\end{aligned}$$

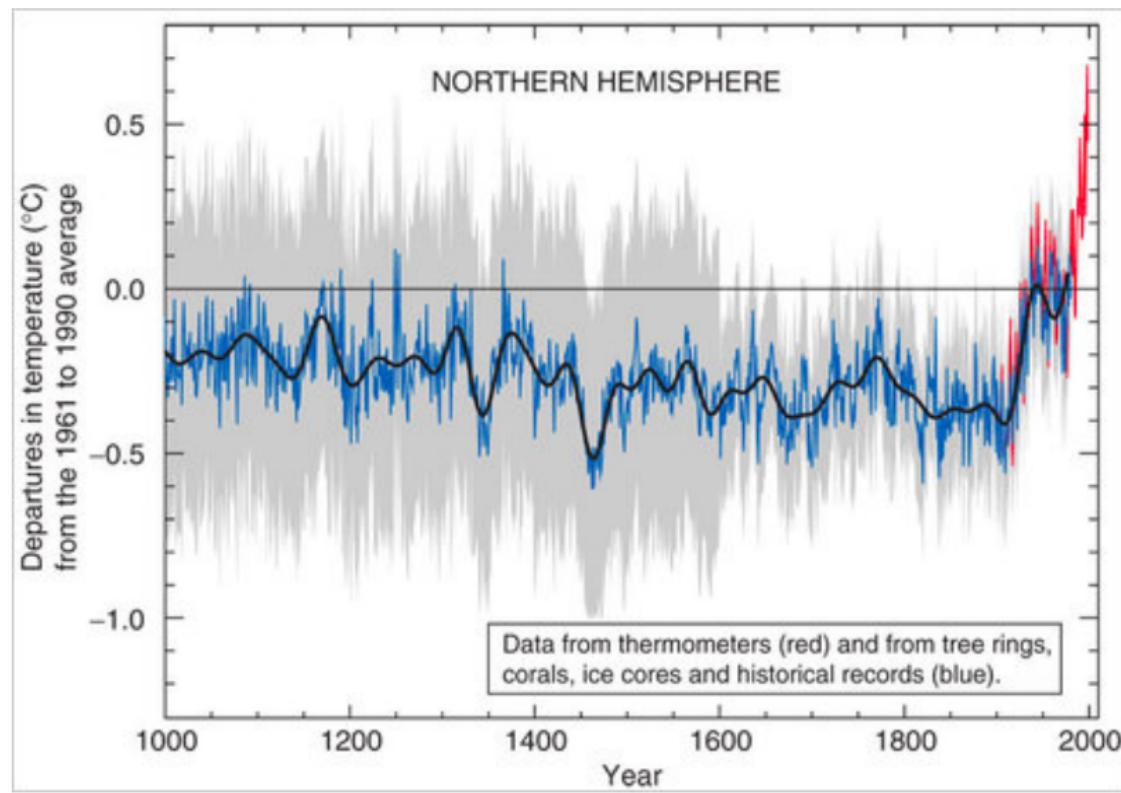
- ▶ The problem gets trickier if f and g above are deterministic models
- ▶ Some quite complex deterministic models have been suggested for pollen/climate. Not many for other proxies
- ▶ Quite a few simple climate models that might work over the palaeoclimate period, e.g. Saltzman and Maasch, 1991:

$$dX_{(1)} = - (X_{(1)} + X_{(2)} + vX_{(3)} + F(\gamma_P, \gamma_C, \gamma_E)) dt + \sigma_1 dW_{(1)}$$

$$dX_{(2)} = (rX_{(2)} - pX_{(3)} - sX_{(2)}^2 - X_{(2)}^3) dt + \sigma_2 dW_{(2)}$$

$$dX_{(3)} = -q(X_{(1)} + X_{(3)}) dt + \sigma_3 dW_{(3)}$$

Back to the future: can we do better than this?



Summary

- ▶ A Bayesian model with an improved forward model and richer climate process for multiple sites and proxies is the ultimate research goal
- ▶ We need help with Bayesian computation for large multivariate non-linear non-Gaussian state space models
- ▶ We need help with combining deterministic/stochastic elements in forward models and climate models
- ▶ We must do better than the Hockey Stick!