

# Bayesian Modeling and Computation in Complex Geophysical Problems

---

Mark Berliner  
Ohio State University

NCAR Workshop: Petascale Computing 5-7 May, 2008

## Outline

- Bayesian Modeling: Selected Features
- Computation: Monte Carlo (MC)
  1. Basic MC
  2. Importance Sampling, Particle Filters
  3. Markov Chain MC
- Examples of Multiscale Models
- Discussion

# Bayesian Modeling: Selected Features

---

- **Bayesian Analysis:** treating uncertainty and knowledge
  - Combine observations & other information sources formally
  - Uncertainty management is paramount
  - Inputs and outputs are probability distributions
- **Mechanism:** probability theory
- **Challenges:** (I) formulation of models; (II) computation.
- **Two general arenas**
  1. **Stochastic Dynamic Modeling:** developing probability models for a complex system (within & across space-time scales; coarse graining; stochastic parameters & parameterizations)
  2. **Forecasting:** learning about and predicting an unobserved trajectory of a dynamical system (NWP; data assimilation)

# Modeling device: Bayesian Hierarchical Models (BHM)

---

- HM: Sequences of conditional probability models

$$p(\mathbf{x}, \mathbf{y}, \mathbf{z}) = p(\mathbf{x} \mid \mathbf{y}, \mathbf{z}) p(\mathbf{y} \mid \mathbf{z}) p(\mathbf{z})$$

- Skeleton BHM, Observations  $\mathbf{y}$ ; Processes  $\mathbf{x}$ ; Parameters  $\theta$

1. Data Model  $q(\mathbf{y} \mid \mathbf{x}, \theta)$

2. Prior Process Model  $p(\mathbf{x} \mid \theta)$

3. Prior Parameter Model  $p(\theta)$

- Bayes' Theorem gives Posterior Distribution:

$$\begin{aligned} p(\mathbf{x}, \theta \mid \mathbf{y}) &\propto q(\mathbf{y} \mid \mathbf{x}, \theta) p(\mathbf{x} \mid \theta) p(\theta) \\ &= q(\mathbf{y} \mid \mathbf{x}, \theta) p(\mathbf{x} \mid \theta) p(\theta) / q(\mathbf{y}) \end{aligned}$$

$$\text{where } q(\mathbf{y}) = \int q(\mathbf{y} \mid \mathbf{x}, \theta) p(\mathbf{x} \mid \theta) p(\theta) d\mathbf{x} d\theta$$

# What Does This Buy Us?

---

- Combining information: “Physical-statistical modeling” (Berliner, JGR, 2003)
- Quantifying and dealing with uncertainty!!
- SGS parameterization? (ECMWF & “stoch-physics”)

A. Operational impact of chaos: treat things as random.

B. From a deterministic physical model,  $\mathcal{D}(\mathbf{x}, \boldsymbol{\theta}_m) = 0$   
to a stochastic model  $p(\mathbf{X} | \boldsymbol{\theta})$

- “Approximate physics ( $\mathcal{D}$ ) applied approximately (discretize  $\mathcal{D}$ ) and unsurely ( $\boldsymbol{\theta}$ ; forcings unknown)”
- (Berliner, Milliff, Wikle, 2003, JGR) Air-sea interaction:

$$\left(\nabla^2 - \frac{1}{r^2}\right) \frac{\partial \psi}{\partial t} = -\mathbf{J}(\psi, \nabla^2 \psi) - \beta \frac{\partial \psi}{\partial x} + \frac{1}{\rho H} \text{curl}_z \tau(\mathbf{W}) - \gamma \nabla^2 \psi + a_h \nabla^4 \psi.$$

I see  $p(\psi_{t+1} | \psi_t, \boldsymbol{\theta}, \text{winds, boundary \& initial con.})$

## C. Parameterization: Physical variables $\mathbf{X}, \mathbf{Z} = (\mathbf{Z}_r, \mathbf{Z}_u)$

- Discrete Time Physical Model:

- $\mathbf{X}_{t+1} = \mathbf{h}(\mathbf{X}_t, \mathbf{Z}_{t+1})$

- $\mathbf{Z}_{t+1} = \mathbf{g}(\mathbf{X}_t, \mathbf{Z}_t)$

- Numerical, parameterized model

- $\mathbf{x}_{t+1} = \tilde{\mathbf{h}}(\mathbf{x}_t, \mathbf{z}_{r,t+1}, \mathbf{z}_{u,t+1})$  and  $\mathbf{z}_{u,t+1} \approx \mathbf{F}(\mathbf{x}_t, \mathbf{z}_{r,t+1}, \boldsymbol{\theta})$  give

- $\mathbf{x}_{t+1} = \mathbf{G}(\mathbf{x}_t, \mathbf{z}_{r,t+1}, \mathbf{F}(\mathbf{x}_t, \mathbf{z}_{r,t+1}, \boldsymbol{\theta}))$

### 1. Stochastic-Bayesian Parameterization

$$\mathbf{p}(\mathbf{x}_{t+1} \mid \mathbf{x}_t, \mathbf{z}_{r,t+1}, \boldsymbol{\theta}) = \int \mathbf{p}(\mathbf{x}_{t+1} \mid \mathbf{x}_t, \mathbf{z}_{u,t+1}, \mathbf{z}_{r,t+1}, \boldsymbol{\theta}) \mathbf{p}(\mathbf{z}_{u,t+1} \mid \mathbf{x}_t, \mathbf{z}_{r,t+1}, \boldsymbol{\theta}) d\mathbf{z}_{u,t+1}$$

### 2. On-the-fly Stochastic-Bayesian Parameterization

$$\mathbf{p}(\mathbf{x}_{t+1} \mid \mathbf{x}_t, \mathbf{z}_{r,t+1}, \boldsymbol{\theta}, \mathbf{Y}) = \int \mathbf{p}(\mathbf{x}_{t+1} \mid \mathbf{x}_t, \mathbf{z}_{u,t+1}, \mathbf{z}_{r,t+1}, \boldsymbol{\theta}, \mathbf{Y}) \mathbf{p}(\mathbf{z}_{u,t+1} \mid \mathbf{x}_t, \mathbf{z}_{r,t+1}, \boldsymbol{\theta}, \mathbf{Y}) d\mathbf{z}_{u,t+1}$$

- Both Bayesian parameterizations are built using observations, model explorations, etc.

# Bayesian Computation and Monte Carlo

---

- Bayes' Theorem gives Posterior Distribution:

$$p(\mathbf{x}, \boldsymbol{\theta} | \mathbf{y}) = q(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta})p(\mathbf{x} | \boldsymbol{\theta})p(\boldsymbol{\theta})/q(\mathbf{y})$$

$$\text{where } q(\mathbf{y}) = \int q(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta})p(\mathbf{x} | \boldsymbol{\theta})p(\boldsymbol{\theta})d\mathbf{x}d\boldsymbol{\theta}$$

- If  $q(\mathbf{y})$  is intractable, turn to Monte Carlo.

## 1. Monte Carlo (MC)

- Sample or ensemble  $\mathbf{x}^1, \dots, \mathbf{x}^M$  from  $p(\mathbf{x} | \mathbf{y})$   
(Suppress  $\boldsymbol{\theta}$ )
- Estimate expectations: (notation:  $E(\ )$  same as  $\langle \rangle$ )

$$E(\mathbf{h}(\mathbf{X}) | \mathbf{y}) = \int \mathbf{h}(\mathbf{x})p(\mathbf{x} | \mathbf{y})d\mathbf{x} \quad \text{by} \quad \hat{E}(\mathbf{h} | \mathbf{y}) = \frac{1}{M} \sum \mathbf{h}(\mathbf{x}^i)$$

- That is, approximate  $p(\mathbf{x} | \mathbf{y})$  by discrete, uniform distribution on the sample:  $\widehat{\Pr}(\mathbf{X} = \mathbf{x}^i) = 1/M$

## 2. Importance Sampling (ISM)

- Direct sampling from  $p(\mathbf{x} | \mathbf{y})$  very hard or not possible
- Sample  $\mathbf{x}^1, \dots, \mathbf{x}^M$  from  $g$
- Estimate

$$\mathbb{E}(h(\mathbf{X}) | \mathbf{y}) = \int h(\mathbf{x}) \frac{p(\mathbf{x} | \mathbf{y})}{g(\mathbf{x})} g(\mathbf{x}) d\mathbf{x} \quad \text{by} \quad \frac{1}{M} \sum h(\mathbf{x}^i) \frac{p(\mathbf{x}^i | \mathbf{y})}{g(\mathbf{x}^i)}$$

- Usual alternative:

- Define normalized ISMC weights  $\alpha_i = w(\mathbf{x}^i) / \sum w(\mathbf{x}^j)$  where  $w(\mathbf{x}^i) = p(\mathbf{x}^i | \mathbf{y}) / g(\mathbf{x}^i)$
- Estimation:  $\hat{\mathbb{E}}(h | \mathbf{y}) = \sum \alpha_i h(\mathbf{x}^i)$
- Approximate  $p(\mathbf{x} | \mathbf{y})$  by discrete distribution  $\{\mathbf{x}^i, \alpha_i : i = 1, \dots, M\}$ :  
 $\widehat{\Pr}(\mathbf{X} = \mathbf{x}^i) = \alpha_i$
- Key: the normalizer  $q(\mathbf{y})$  of  $p(\mathbf{x} | \mathbf{y})$  cancels in the  $\alpha$ 's so we only need  $p(\mathbf{x} | \mathbf{y})$  up to proportionality.

## Notes on ISMC

- **Particle Filter: Evolve or generate  $\mathbf{x}_t^i$  over time.**
  - Sample  $\{\mathbf{x}_{t-1}^i, \alpha_{i,t-1} : i = 1, \dots, M\}$  representing  $p(\mathbf{x}_{t-1} | \mathbf{y}_{t-1})$
  - Generate  $\mathbf{x}_t^i \sim p(\mathbf{x}_t | \mathbf{x}_{t-1}^i)$   
 $\{\mathbf{x}_t^i, \alpha_{i,t-1} : i = 1, \dots, M\}$  represents  $p(\mathbf{x}_t | \mathbf{y}_{t-1})$   
(Forecast Step)
  - Bayes' Theorem converts to  $\{\mathbf{x}_t^i, \alpha_{i,t} : i = 1, \dots, m\}$  representing  $p(\mathbf{x}_t | \mathbf{y}_t)$  where
$$\alpha_{i,t} \propto q(\mathbf{y}_t | \mathbf{x}_t^i) \alpha_{i,t-1}$$
(Analysis Step)
- What we can do with an ensemble depends on how it was made.
- In high dimensions  $\alpha$ 's are poorly behaved:  
They concentrate on a few (or one!) ensemble members



### 3. Markov Chain Monte Carlo (MCMC)

- Finding normalizer  $q(y)$  vrs finding “partition function” in Statistical Mechanics
- MCMC: Develop a stationary (ergodic) Markov chain with limiting distribution coinciding with the target posterior  $p(x | y)$ .
  - After a “burn-in” (like “spin-up”) period, realizations from the chain form an ensemble from  $p(x | y)$  (approximately).
  - Ensemble members are dependent, but MC estimation works

## Metropolis-Hastings

- State of chain at iterate  $i$  :  $\mathbf{x}^i$ 
  - generate  $\tilde{\mathbf{x}}$  from some proposal distribution  $g(\tilde{\mathbf{x}} | \mathbf{x}^i)$
  - generate independent  $U = \text{Uniform}(0,1)$  RV
  - set  $\mathbf{x}^{i+1} = \tilde{\mathbf{x}}$  if
$$U < \frac{p(\tilde{\mathbf{x}} | \mathbf{y}) g(\mathbf{x}^i | \tilde{\mathbf{x}})}{p(\mathbf{x}^i | \mathbf{y}) g(\tilde{\mathbf{x}} | \mathbf{x}^i)}$$
  - set  $\mathbf{x}^{i+1} = \mathbf{x}^i$ , otherwise.
- Key: Again, normalizer  $q(\mathbf{y})$  cancels.

## Gibbs Sampler

- $\mathbf{x}$  is a  $K$ -dimensional vector,  $(\mathbf{x}_1, \dots, \mathbf{x}_K)$
- Derive “full conditionals”  $p(\mathbf{x}_k \mid \mathbf{x}_1, \dots, \mathbf{x}_{k-1}, \mathbf{x}_{k+1}, \dots, \mathbf{x}_K \text{ (and } \mathbf{y}))$ 
  - state of chain at iterate  $i$ :  $(\mathbf{x}_1^i, \dots, \mathbf{x}_K^i)$
  - generate  $\mathbf{x}_1^{i+1}$  from  $p(\mathbf{x}_1 \mid \mathbf{x}_2^i, \mathbf{x}_3^i, \dots, \mathbf{x}_K^i)$
  - generate  $\mathbf{x}_2^{i+1}$  from  $p(\mathbf{x}_2 \mid \mathbf{x}_1^{i+1}, \mathbf{x}_3^i, \dots, \mathbf{x}_K^i)$
  - ⋮
  - generate  $\mathbf{x}_K^{i+1}$  from  $p(\mathbf{x}_K \mid \mathbf{x}_1^{i+1}, \dots, \mathbf{x}_{K-1}^{i+1})$

## Others

- Metropolis-within-Gibbs: replace intractable full conditionals by Metropolis steps.
- Using stochastic differential equation

$$du(t) = b(u)dt + \sigma(u)dW(t)$$

where  $\{W(t) : t \geq 0\}$

- (Theory & assumptions)  $U(t)$  has a density function  $p(u, t)$
- $p(u, t)$  is solution Fokker-Planck Equation
- Stationary solutions:

$$0.50 \frac{\partial^2}{\partial u^2}(\sigma^2 p) = \frac{\partial}{\partial u}(b p)$$

- Pick  $b$  and  $\sigma$  so that  $p$  is the target posterior
- ISMC-MCMC

## Key Points

- Monitoring convergence
- Output Analysis: Using output to summarize the target posterior.
- Tensions:
  - Multiple runs vrs one long run  
Wasted burnin periods vrs “mixing”
  - Output from a run are dependent

$$\text{Var}(\bar{\mathbf{x}}^i) = \frac{\mathbf{v}^2}{M}(\mathbf{1} + \sum \mathbf{c}(\mathbf{i})\rho_i)$$

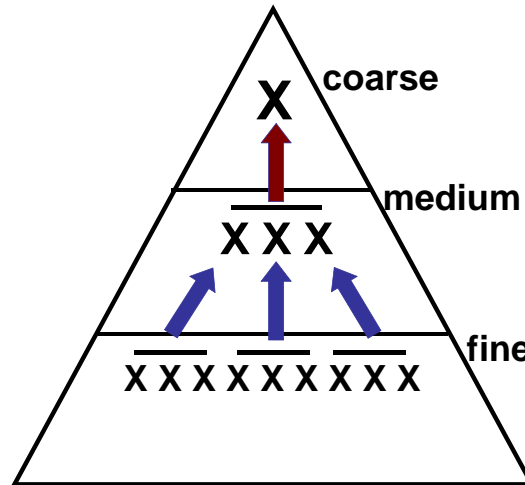
- Embarassingly Parallel? Obvious: Multiple runs, but?

# Two Notions of Multiscale Modeling

---

## 1. Space-Time Filtering

- “Hierarchical” process prior:  $p(\vec{X}_c | \vec{X}_m) p(\vec{X}_m | \vec{X}_f) p(\vec{X}_f)$
- Up-down scaling:  $p(\vec{X}_c, \vec{X}_f | \vec{X}_m) p(\vec{X}_m)$
- Terra incognita:  $p(\vec{X}_m | \vec{X}_f, \vec{X}_c) p(\vec{X}_f | \vec{X}_c)$

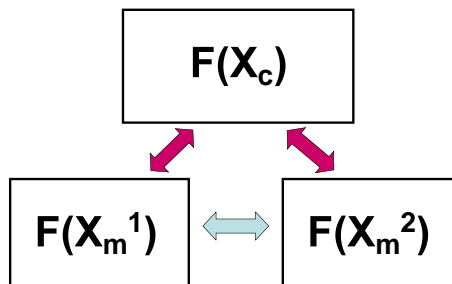


## Example

- **Data Model:**  $q(\vec{Y}_c | \vec{X}_c) q(\vec{Y}_m^1 | \vec{X}_m^1) q(\vec{Y}_m^2 | \vec{X}_m^2)$
- **Process Prior:**  $p(\vec{X}_c | \vec{X}_m^1, \vec{X}_m^2) p(\vec{X}_m^2 | \vec{X}_m^1) p(\vec{X}_m^1)$

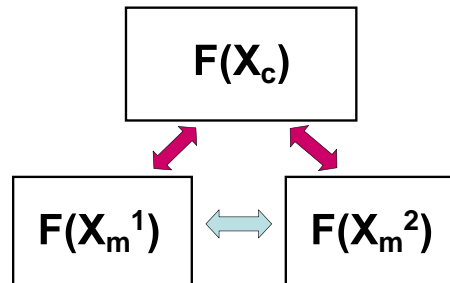
### Full Conditionals for Gibbs Sampler:

- $F(\vec{X}_c | \text{rest}) \propto q(\vec{Y}_c | \vec{X}_c) p(\vec{X}_c | \vec{X}_m^1, \vec{X}_m^2)$
- $F(\vec{X}_m^2 | \text{rest}) \propto q(\vec{Y}_m^2 | \vec{X}_m^2) p(\vec{X}_c | \vec{X}_m^1, \vec{X}_m^2) p(\vec{X}_m^2 | \vec{X}_m^1)$
- $F(\vec{X}_m^1 | \text{rest}) \propto q(\vec{Y}_m^1 | \vec{X}_m^1) p(\vec{X}_c | \vec{X}_m^1, \vec{X}_m^2) p(\vec{X}_m^2 | \vec{X}_m^1) p(\vec{X}_m^1)$
- **Note how all levels intertwine: challenge to parallel code**



## Example Cont'd: Potentials for parallel codes

- (1) Run bottom nodes holding  $\vec{X}_c$  fixed  
(i.e., we don't have to update every variable every time, though not doing so may slow convergence/)
- Master swaps across scales occasionally.
  - Many scales: Management system

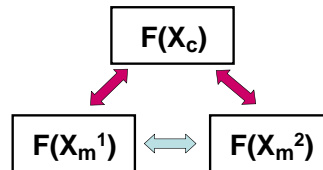




## Example Cont'd: Potentials for parallel codes

### (2) Partial Conditionals for Gibbs Sampler with ISMC:

- $F(\vec{X}_c \mid \text{rest}) \propto q(\vec{Y}_c \mid \vec{X}_c) p(\vec{X}_c \mid \vec{X}_m^1, \vec{X}_m^2)$
- $F_p(\vec{X}_m^2 \mid \text{rest}) \propto q(\vec{Y}_m^2 \mid \vec{X}_m^2) p(\vec{X}_m^2 \mid \vec{X}_m^1) \{p(\vec{X}_c \mid \vec{X}_m^1, \vec{X}_m^2)\}$
- $F_p(\vec{X}_m^1 \mid \text{rest}) \propto q(\vec{Y}_m^1 \mid \vec{X}_m^1) p(\vec{X}_m^1) \{p(\vec{X}_c \mid \vec{X}_m^1, \vec{X}_m^2) p(\vec{X}_m^2 \mid \vec{X}_m^1)\}$ 
  - Ignoring terms in brackets if results are simple
  - But, those terms form required IS weights
  - Speed versus memory



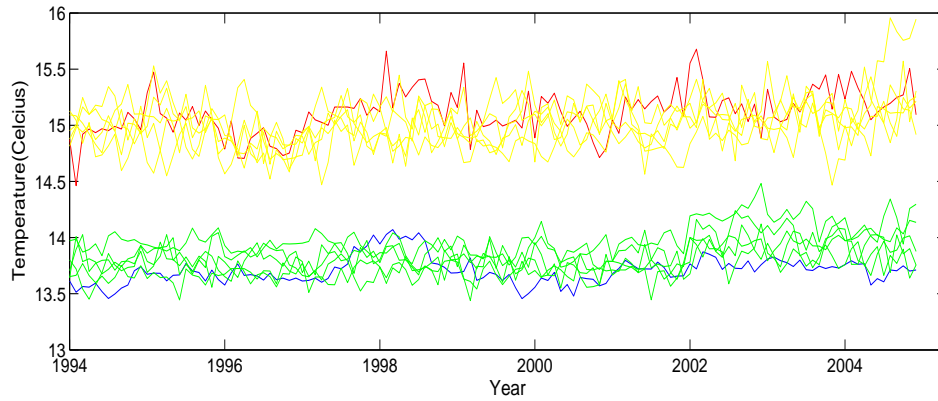
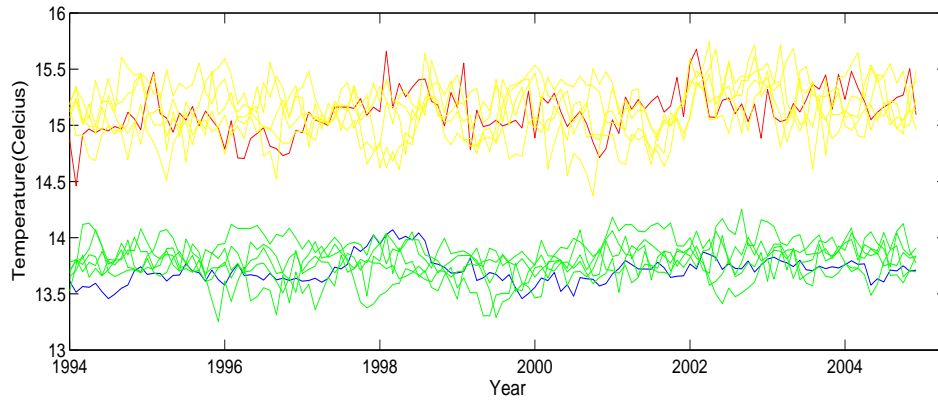
## 2. Parameterization of Scales

- Build or “parameterize” scales into dynamic model for  $\mathbf{X}$   
Example (Berliner & Kim, 2008, J Clim)
- $\mathbf{X}$ : monthly surface temperatures
- Time series models (AR) with time varying parameters

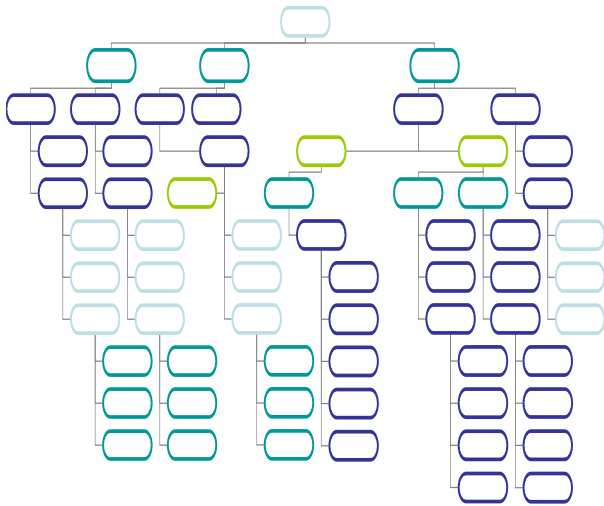
$$\mathbf{X}_t = \boldsymbol{\mu}_{i(t)} + \boldsymbol{\beta}_{j(t)}(\mathbf{X}_{t-1} - \boldsymbol{\mu}_{i(t-1)}) + \mathbf{e}_{(t)}$$

- $\boldsymbol{\mu}_{i(t)}$  slowly vary (climate scale);  $\boldsymbol{\beta}_{j(t)}$  vary moderately (another climate scale);  $\mathbf{e}_t$  vary quickly (“weather”), but their variances slowly vary (climate scale):
  - $\boldsymbol{\mu}_i = \mathbf{a} + \mathbf{b} \text{CO}_{2i} + \text{noise}$
  - $\boldsymbol{\beta}_j = \mathbf{c} + \mathbf{d} \text{SOI}_j + \text{noise}$

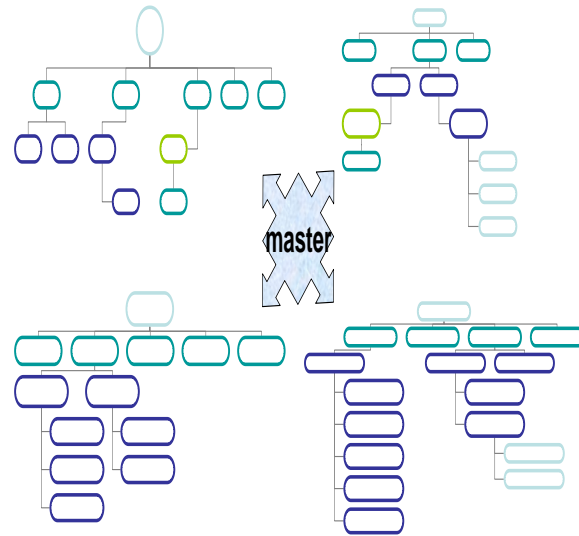
- **Computational challenge: Model selection**  
With what rates should the  $\mu_{i(t)}$ ,  $\beta_{j(t)}$ ,  
and variances of the  $e_t$  evolve?  
(1000's of combinations)
- **Decadal Prediction**
  - Build model using observations up to 1994
  - Forecast for the following 10 years using a stochastic model for SOI
  - Next Graphic: show NH and SH observed temp's and ensembles from our predictive distributions  
(First panel:  $\mu_{i(t)}$  varying every 8 years,  
 $\beta_{j(t)}$  varying every 2 years  
second panel:  $\mu_{i(t)}$  varying every 8 years,  
 $\beta_{j(t)}$  varying every 4 years)



## 1. Bayesian Networks



## 2. Competing Networks



# Discussion

---

- Joe's talk