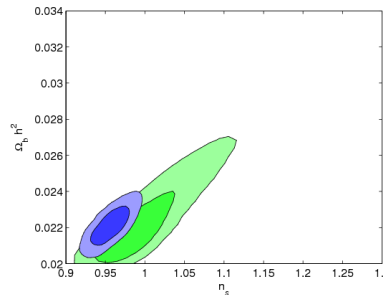


# Computing posteriors

- For 2 parameters, a grid is usually possible
  - Marginalise by numerically integrating along each axis of the grid



- For  $\gg 2$  parameters it is not feasible to have a grid (e.g. 10 points in each parameter direction, 12 parameters =  $10^{12}$  likelihood evaluations)
- We can instead *sample* the distribution

# Sampling

- Draw samples with a probability proportional to the *target distribution* (likelihood, or the posterior, or something else)
- Approximate target by a sum of Dirac delta functions:

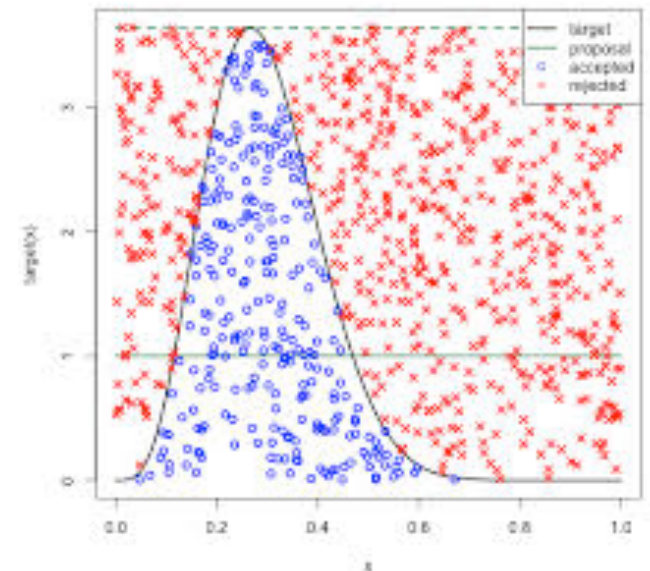
$$p(\theta) \simeq \frac{1}{N} \sum_{\text{samples } i} \delta^D(\theta - \theta_i)$$

- Mean of any function is given approximately by

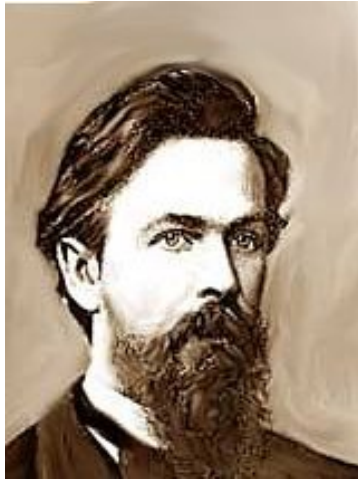
$$\langle f \rangle = \int f(\theta) p(\theta) d\theta \simeq \frac{1}{N} \sum_{\text{samples } i} f(\theta_i)$$

# Advantages of sampling

- Points may be concentrated where the target is high
- Converges as number of samples  $\rightarrow \infty$
- e.g. Rejection sampling:
- Very inefficient in general



Credit: J. Niemi



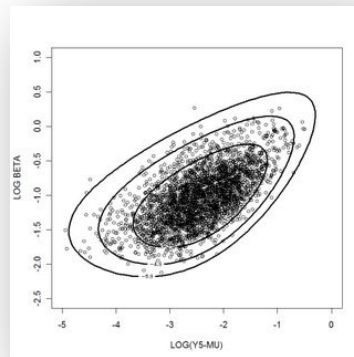
# MCMC

## Markov Chain

## Monte Carlo



Aim of MCMC: generate a set of points in the parameter space whose distribution function is the same as the target density.



MCMC follows a Markov process - i.e. the next sample depends on the present one, but not on previous ones.

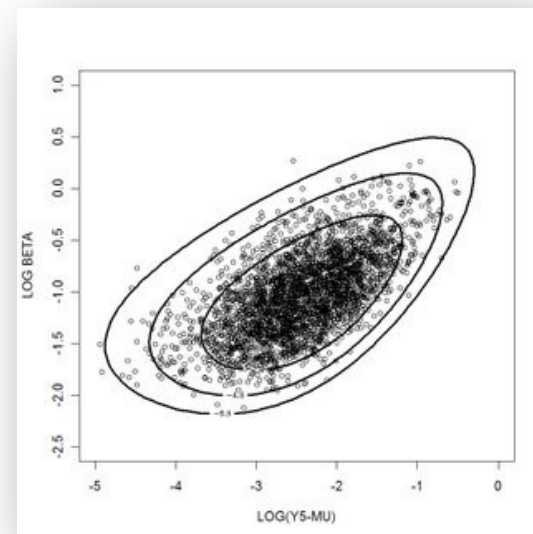
# Target density

The *target density* is approximated by a set of delta functions (you may need to normalise):

$$p(\boldsymbol{\theta}) \simeq \frac{1}{N} \sum_{i=1}^N \delta(\boldsymbol{\theta} - \boldsymbol{\theta}_i)$$

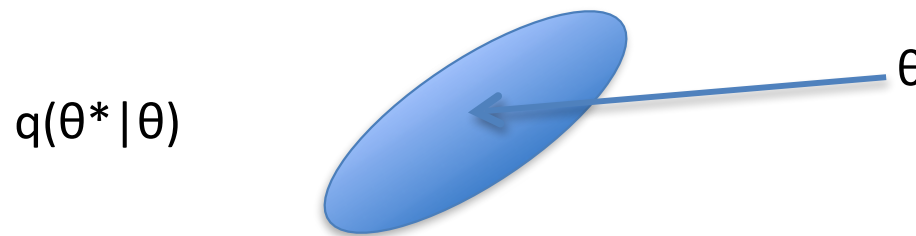
and we can estimate the mean of any function  $f$  by

$$f(\boldsymbol{\theta}) \simeq \frac{1}{N} \sum_{i=1}^N f(\boldsymbol{\theta}_i).$$



# Metropolis-Hastings algorithm

Proposal Distribution



$$p(\text{acceptance}) = \min \left[ 1, \frac{p(\theta^*)q(\theta^* | \theta)}{p(\theta)q(\theta | \theta^*)} \right]$$

Metropolis algorithm (special case):

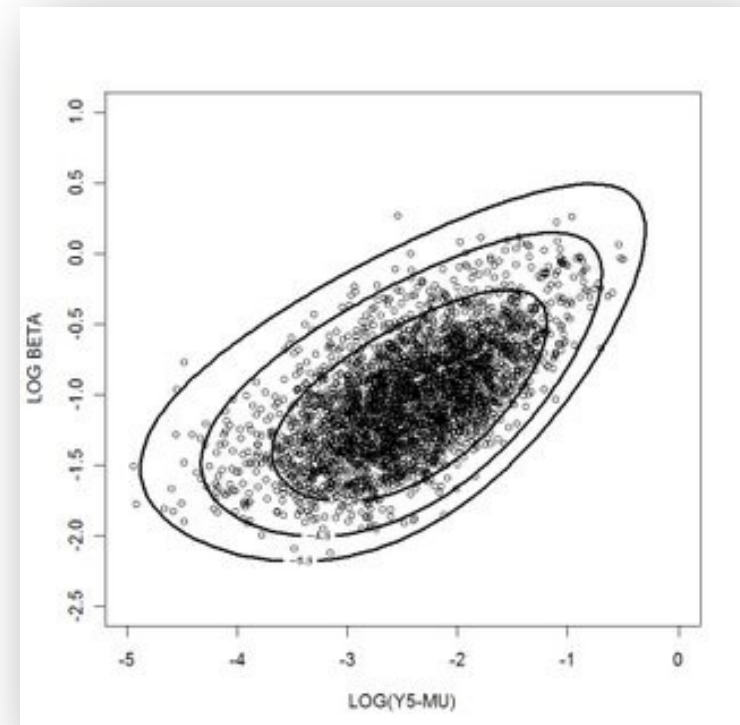
$$\min \left[ 1, \frac{p(\theta^*)}{p(\theta)} \right]$$

# MCMC Algorithm

- Choose a random initial starting point in parameter space, and compute the target density.
- Repeat:
  - Generate a step in parameter space from a proposal distribution, generating a new trial point for the chain.
  - Compute the target density at the new point, and accept it (or not) with the Metropolis-Hastings algorithm.
  - If the point is not accepted, the previous point is repeated in the chain.
- End Repeat:

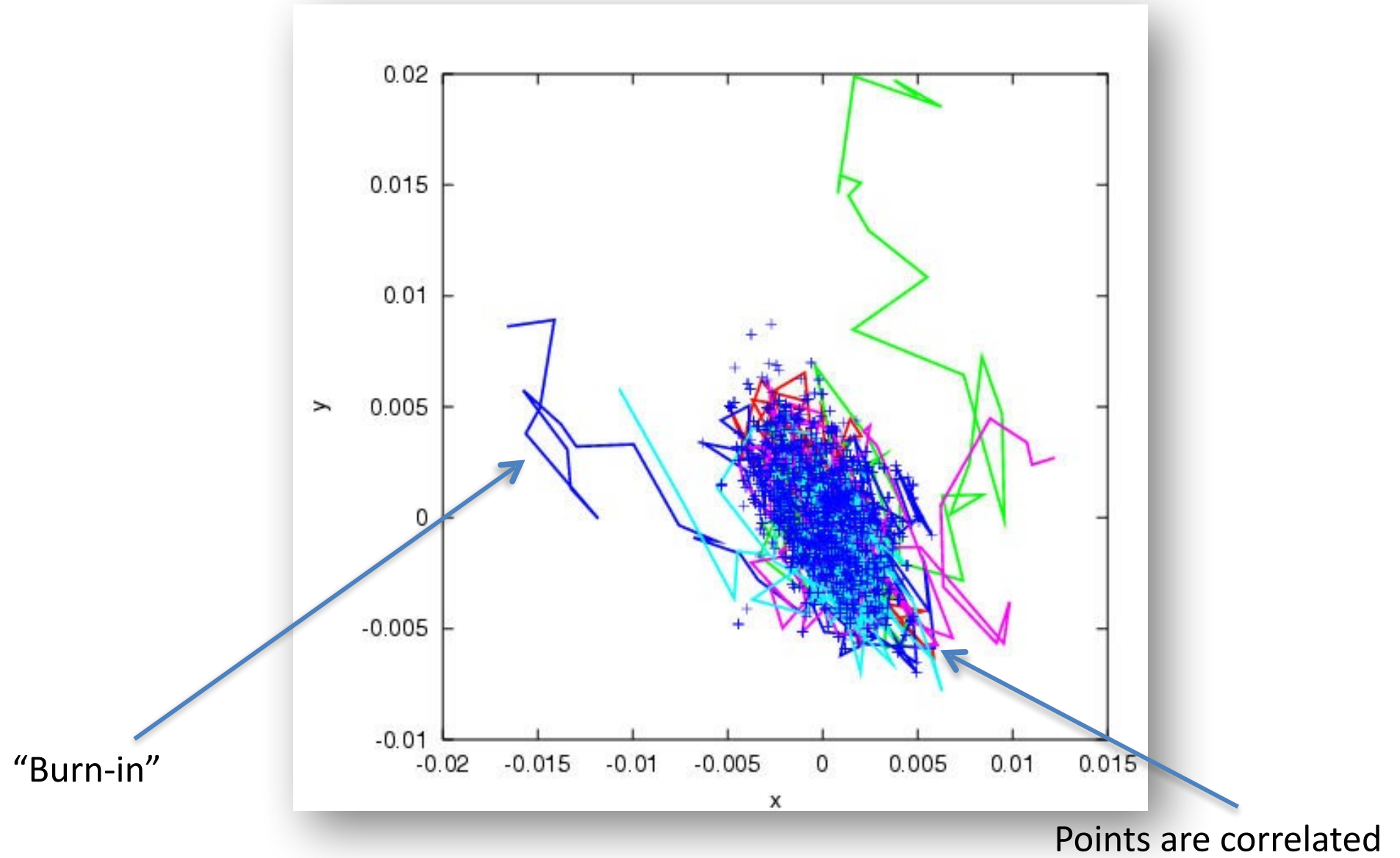
# The proposal distribution

- Too small, and it takes a long time to explore the target
- Too large and almost all trials are rejected
- $q \sim$  'Fisher size' is good.



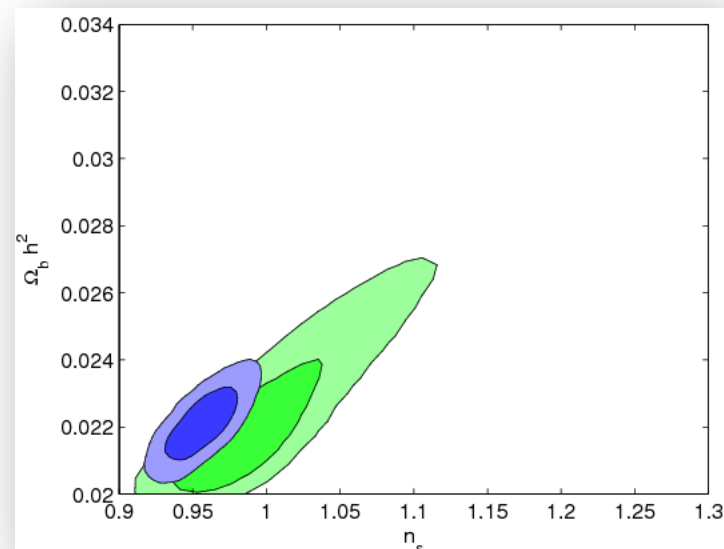
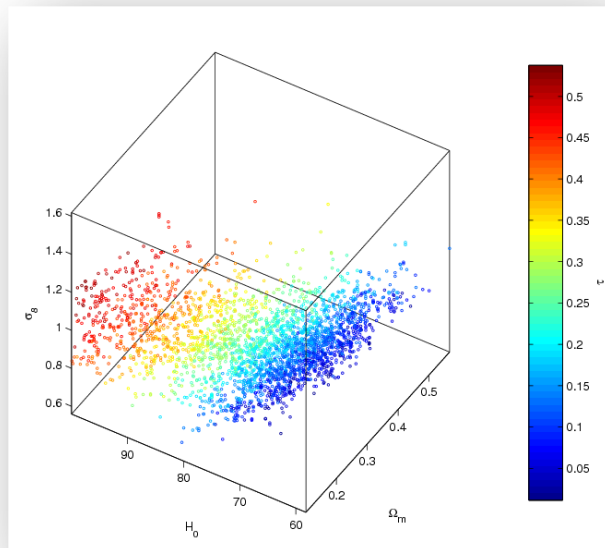


# Burn-in

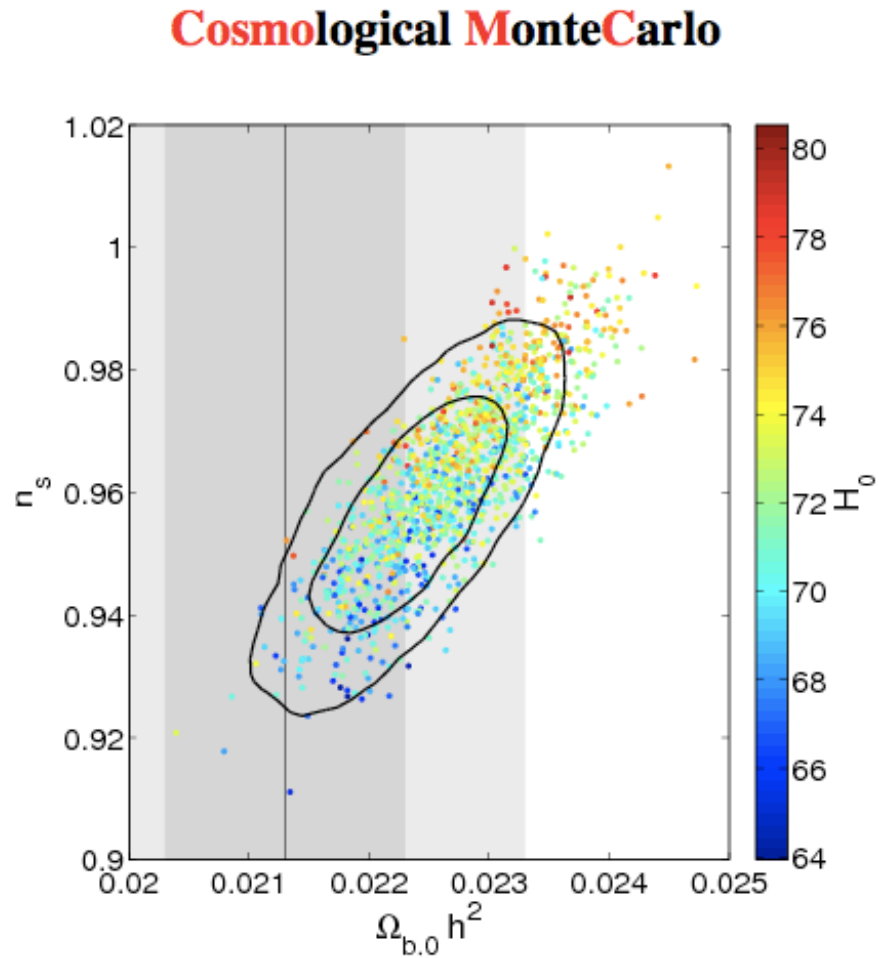


# Marginalisation

- Marginalisation is trivial
  - Each point in the chain is labelled by all the parameters
  - To marginalise, just ignore the labels you don't want



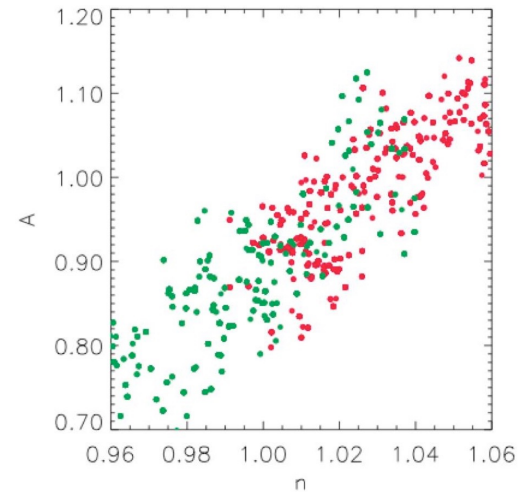
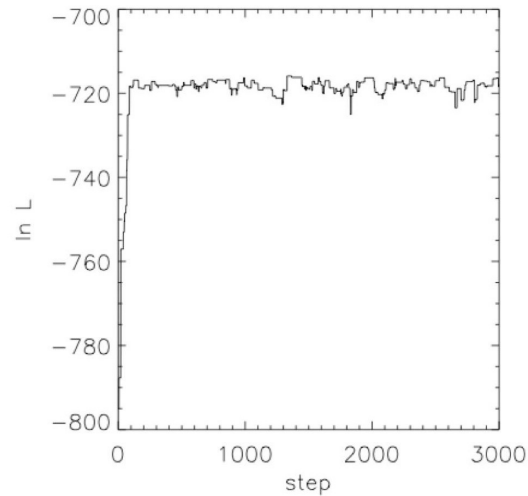
# CosmoMC



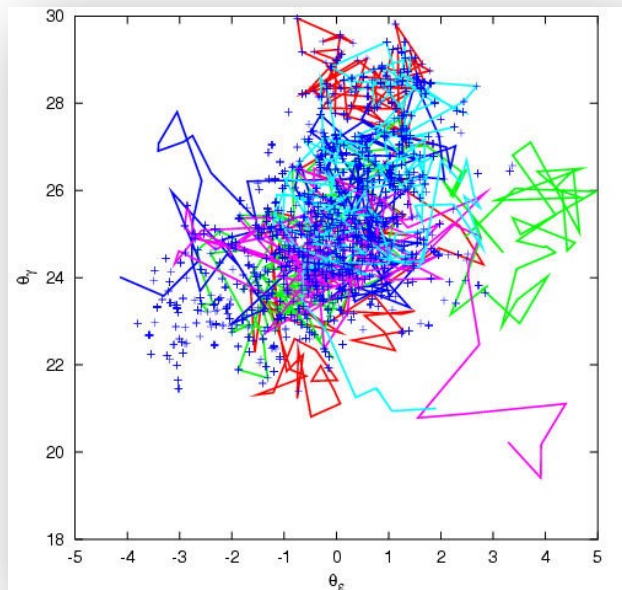
Samples from WMAP 5-yr likelihood combined with deuterium constraint ([0805.0594](#))

<http://cosmologist.info/cosmomc/>

# Convergence



Verde et al 2003



You must use a convergence test.  
Gelman-Rubin test is most common

# Gelman-Rubin test (1992)

- Run  $M$  chains, each of length  $N$ , from different starting points
- Basic idea: compare different estimates of the variance of the points. If the chains have converged, they should agree well

# Gelman-Rubin algorithm

## and Rubin Diagnostics

and Rubin diagnostics (Gelman and Rubin; 1992; Brooks and Gelman; 1997) are based on analyzing multiple simulated MCMC chains by comparing variances within each chain and the variance between chains. Large deviation between these two variances indicates nonconvergence.

Let  $\theta^t$ , where  $t = 1, \dots, n$ , to be the collection of a single Markov chain output. The parameter  $\theta^t$  is the  $t$ th sample of the Markov chain. For notational simplicity, we assume  $\theta^t$  to be single dimensional in this section.

Suppose you have  $M$  parallel MCMC chains that were initialized from various parts of the target distribution. Each chain is of length  $n$  (after discarding the burn-in). The simulations are labeled as  $\theta_m^t$ , where  $t = 1, \dots, n$  and  $m = 1, \dots, M$ . The between-chain variance  $B$  and the within-chain variance  $W$  are calculated as follows:

$$B = \frac{n}{M-1} \sum_{m=1}^M (\bar{\theta}_m - \bar{\theta})^2, \text{ where } \bar{\theta}_m = \frac{1}{n} \sum_{t=1}^n \theta_m^t, \bar{\theta} = \frac{1}{M} \sum_{m=1}^M \bar{\theta}_m$$

$$W = \frac{1}{M} \sum_{m=1}^M s_m^2, \text{ where } s_m^2 = \frac{1}{n-1} \sum_{t=1}^n (\theta_m^t - \bar{\theta}_m)^2$$

The posterior marginal variance,  $\text{var}(\theta | y)$ , is a weighted average of  $W$  and  $B$ . The estimate of the variance is

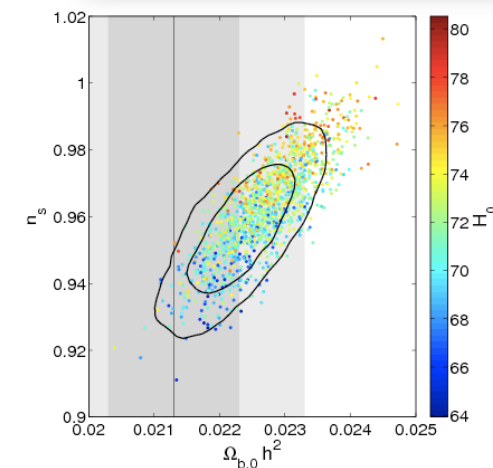
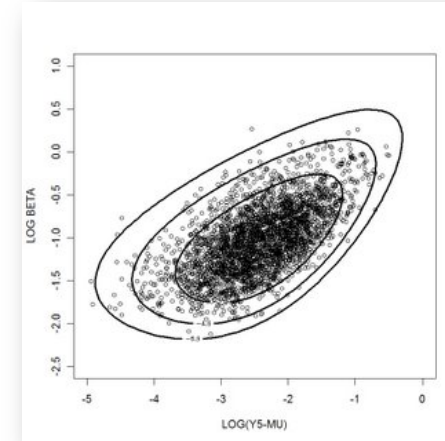
$$\frac{n-1}{n} W + \frac{M+1}{nM} B$$

If the chains have reached the target distribution, this posterior variance estimate should be very close to the within-chain variance  $W$ . Therefore, you would expect the ratio  $\hat{V}/W$  to be close to 1. The square root of this ratio is referred to as the *potential scale reduction factor* (PSRF). A large PSRF indicates that the between-chain variance is substantially greater than the within-chain variance, so that longer simulation is needed. If the PSRF is close to 1, you can conclude that the  $M$  chains have stabilized, and they are likely to have reached the target distribution.

From SAS STAT user guide

# Combining Chains - Importance Sampling

- This is **NOT TRIVIAL**
- For 2 independent experiments, the likelihoods  $L_1$  and  $L_2$  *multiply*
- For small numbers of parameters, chains can be smoothed and multiplied
- For many dimensions, can do **IMPORTANCE SAMPLING**:  
compute  $L_2$  at the positions of the  $L_1$  chain, and give the point a weight  $L_2$ .



Samples from WMAP 5-yr likelihood combined with deuterium constraint ([0805.0594](#))

(Importance sampling draws from a distribution different from the target)

# Gibbs Sampling

- Sampling from a high-dimensional space without symmetries is very hard to do efficiently
- Gibbs sampling can be applied if *conditional* distributions are known, i.e.,

$$p(\theta_i | \text{all other } \theta_j)$$

- Then one can cycle through each parameter in turn
- Can be especially useful for hierarchical models (see later)
- Can sometimes combine with Metropolis-Hastings
- JAGS: <http://mcmc-jags.sourceforge.net/>
- BUGS: [www.mrc-bsu.cam.ac.uk/software/bugs/](http://www.mrc-bsu.cam.ac.uk/software/bugs/)



# Hybrid (Hamiltonian) Monte Carlo

- We would like to increase the acceptance rate to improve efficiency, and explore the target distribution efficiently ('good mixing')
- We have a hard problem in many dimensions. Solution:
- **Make things harder!** add in  $M$  auxiliary variables, one for each parameter in the model.
- Imagine each of the parameters in the problem as a coordinate.
- Target distribution  $\rightarrow$  effective potential
- For each coordinate HMC generates a generalised momentum.
- It then samples from the extended target distribution in  $2M$  dimensions.

# HMC

- HMC explores this  $2M$ -dimensional space by treating the problem as a dynamical system, and evolving the phase space coordinates by solving the dynamical equations.
- Finally, it ignores the momenta (marginalising, as in MCMC), and this gives a sample of the original target distribution.
- May help with decorrelating the points in the chain.
- Invented by particle physicists (Duan et al 1987)

# Theory

- Potential  $U(\theta) = -\ln p(\theta)$
- For each  $\theta_\alpha$ , generate a *momentum*  $\mathbf{u}_\alpha$  (gaussian distributed)
- K.E.  $K = \mathbf{u}^\top \mathbf{u} / 2$
- Define a  $H(\boldsymbol{\theta}, \mathbf{u}) \equiv U(\boldsymbol{\theta}) + K(\mathbf{u})$
- and define an *extended target density*

$$p(\boldsymbol{\theta}, \mathbf{u}) = \exp[-H(\boldsymbol{\theta}, \mathbf{u})]$$

# Magic of HMC

- Evolve as a *dynamical system*

$$\begin{aligned}\dot{\theta}_{\alpha} &= u_{\alpha} \\ \dot{u}_{\alpha} &= -\frac{\partial H}{\partial \theta_{\alpha}}\end{aligned}$$



William Rowan Hamilton

- H remains constant, so extended target density is uniform – *all points get accepted!*
- Also, you can make big jumps – good mixing, if you generate a new  $u$  each time

# Complications

- Evolving the system takes time. Take big steps.
- Need derivatives of  $U = -\ln p$
- We may not have them, and *might approximate*  $U$  (from a short MCMC)

$$U = \frac{1}{2}(\theta - \theta_0)_\alpha C_{\alpha\beta}^{-1}(\theta - \theta_0)_\beta$$

- $H$  is therefore not constant
- Use Metropolis-Hastings. Accept new point with probability

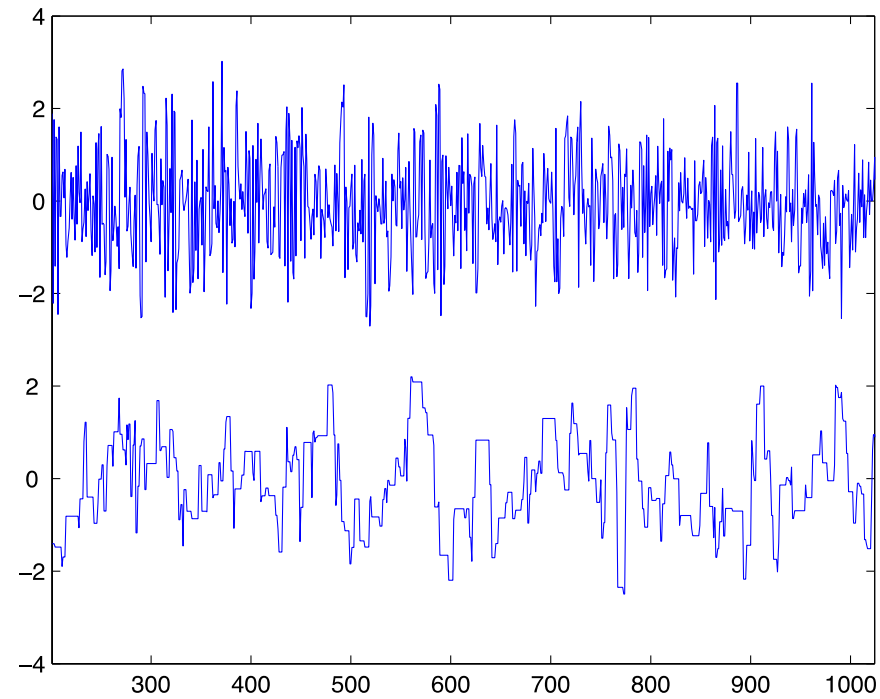
$$\min \{1, \exp [-H(\boldsymbol{\theta}^*, \mathbf{u}^*) + H(\boldsymbol{\theta}, \mathbf{u})]\}$$

# Algorithm

## Hamiltonian Monte Carlo

```
1: initialize  $\mathbf{x}_{(0)}$ 
2: for  $i = 1$  to  $N_{samples}$ 
3:      $\mathbf{u} \sim \mathcal{N}(0, 1)$ 
4:      $(\mathbf{x}_{(0)}^*, \mathbf{u}_{(0)}^*) = (\mathbf{x}_{(i-1)}, \mathbf{u})$ 
5:     for  $j = 1$  to  $N$ 
6:         make a leapfrog move:  $(\mathbf{x}_{(j-1)}^*, \mathbf{u}_{(j-1)}^*) \rightarrow$ 
            $(\mathbf{x}_{(j)}^*, \mathbf{u}_{(j)}^*)$ 
7:     end for
8:      $(\mathbf{x}^*, \mathbf{u}^*) = (\mathbf{x}_{(N)}^*, \mathbf{u}_{(N)}^*)$ 
9:     draw  $\alpha \sim (0, 1)$ 
10:    if  $\alpha < \min\{1, e^{-(H(\mathbf{x}^*, \mathbf{u}^*) - H(\mathbf{x}, \mathbf{u}))}\}$ 
11:         $\mathbf{x}_{(i)} = \mathbf{x}^*$ 
12:    else
13:         $\mathbf{x}_{(i)} = \mathbf{x}_{(i-1)}$ 
14: end for
```

# HMC vs MCMC



Typical speed-ups: factor 4.

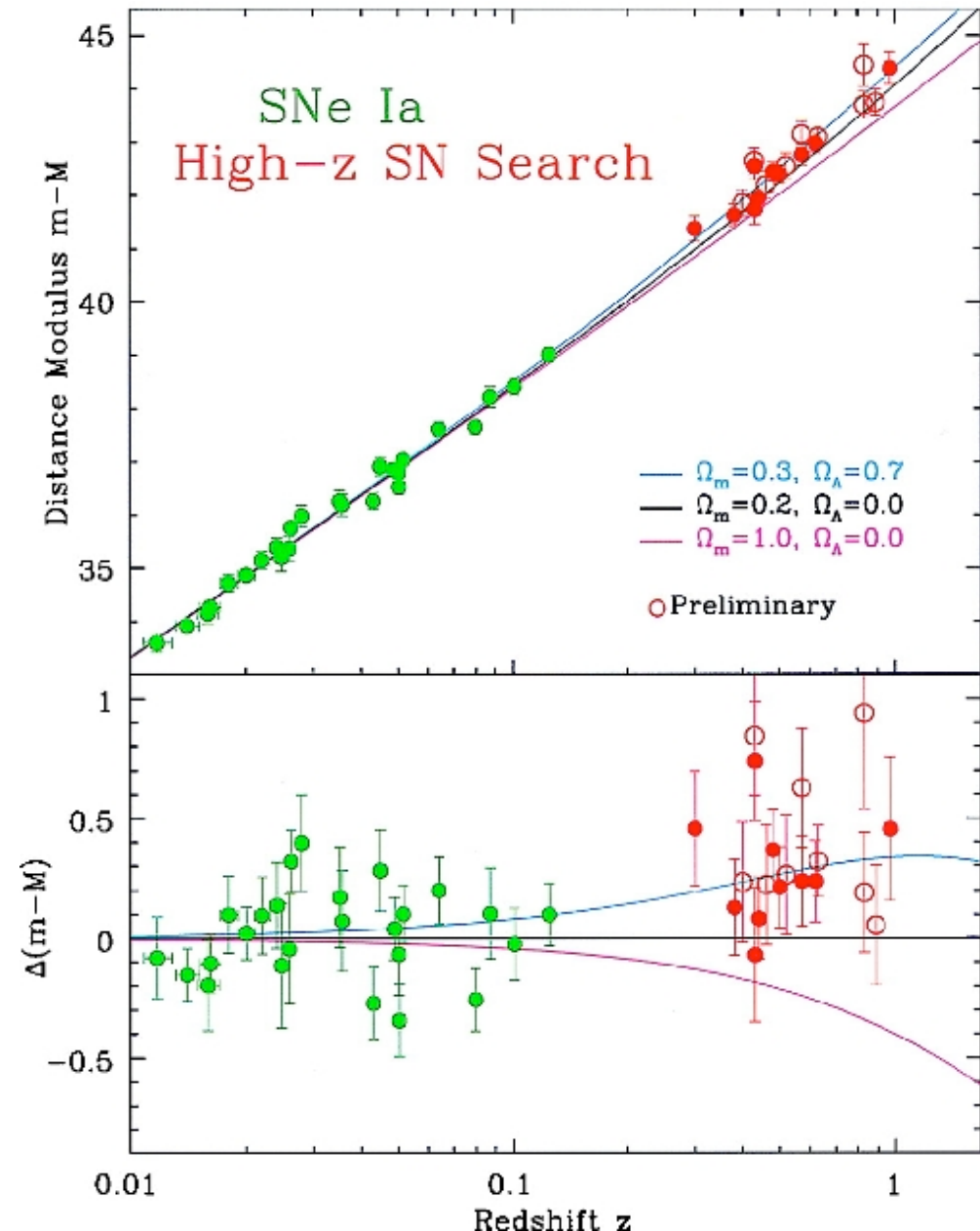
# Supernova distances



- Luminosity Distance depends on  $\Omega_m$  and  $H_0$  (for flat Universe)

$$f = \frac{L}{4\pi D_L^2}$$

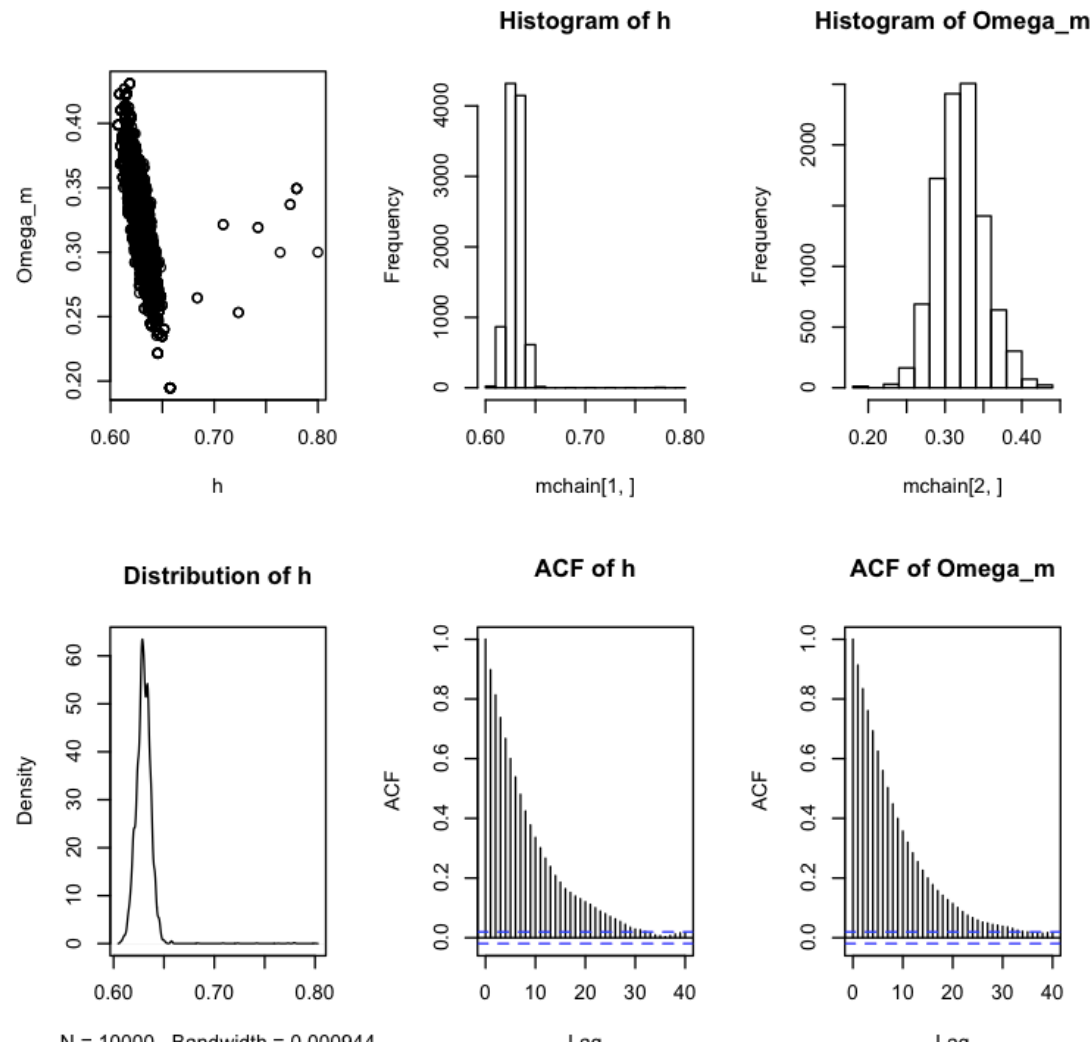
$$D_L = \frac{c(1+z)}{H_0} \int_0^z \frac{dz'}{\sqrt{\Omega_m(1+z')^3 + 1 - \Omega_m}}$$





# MCMC run

Quartz 2 [\*]

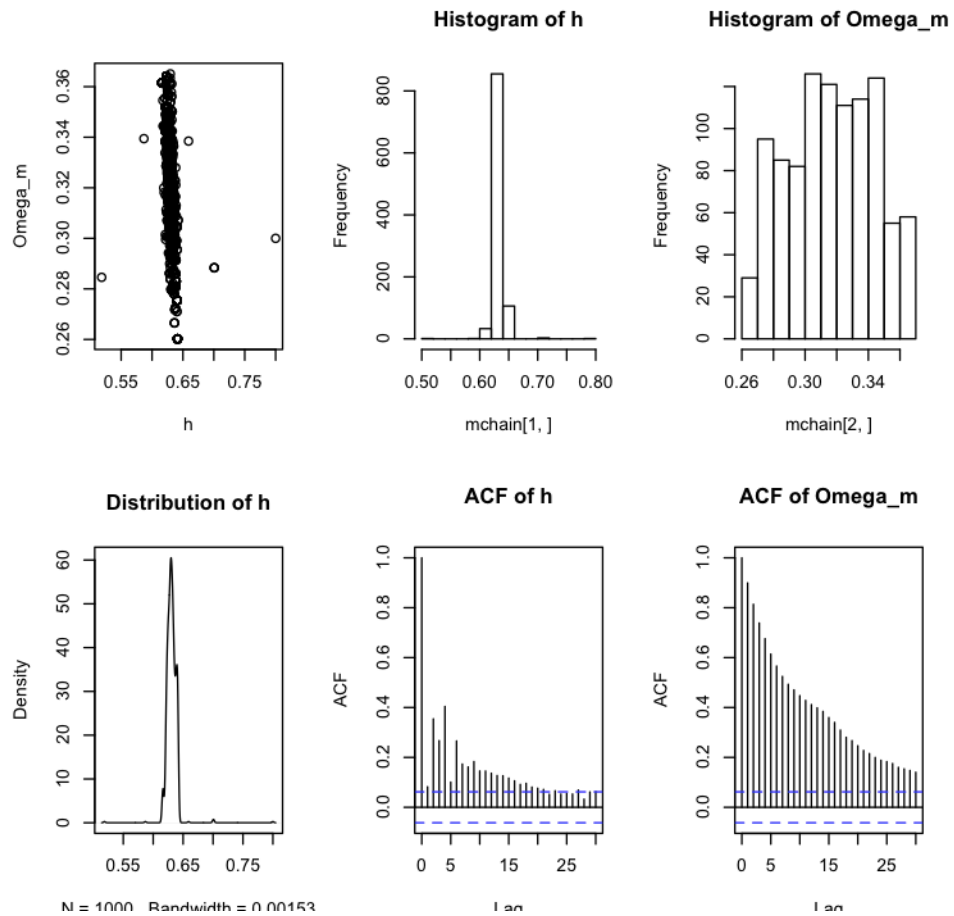


$$C = \begin{pmatrix} 7.1 \times 10^{-5} & -1.9 \times 10^{-4} \\ -1.9 \times 10^{-4} & 1.0 \times 10^{-3} \end{pmatrix}$$

Acceptance  
rate 0.15

# HMC

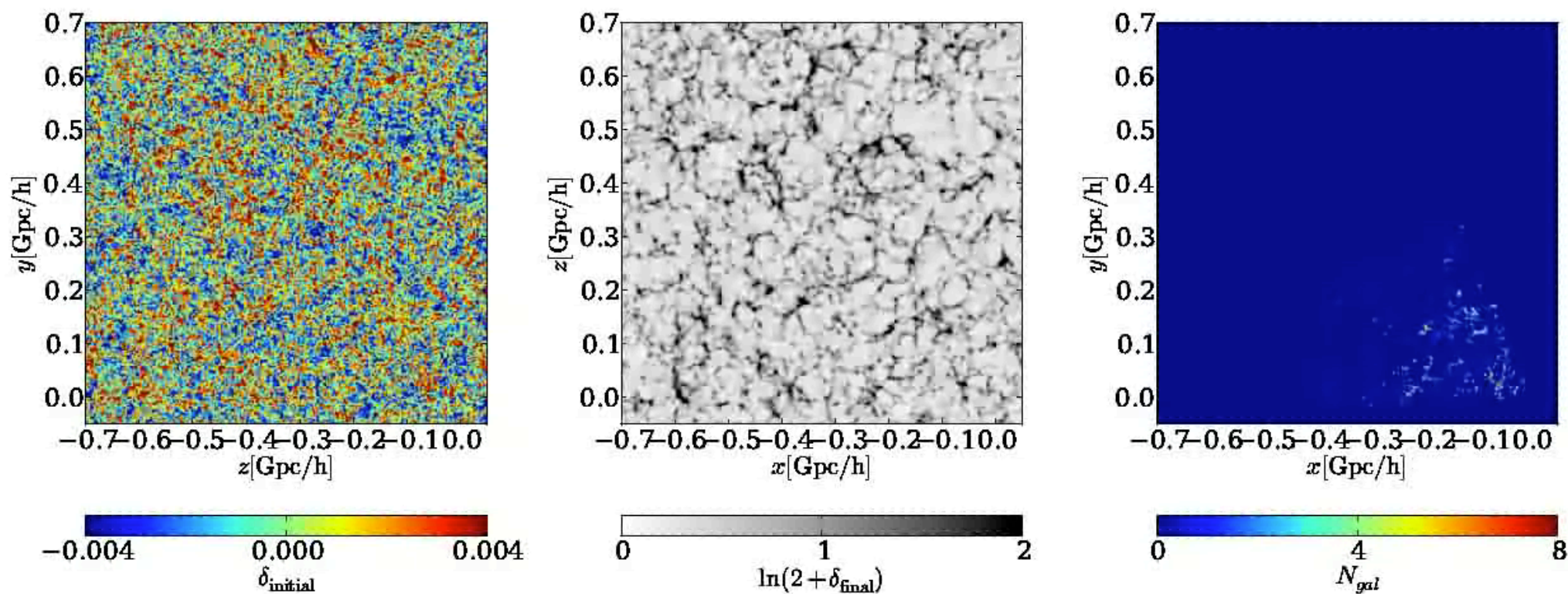
Quartz 2 [\*]



Acceptance rate 0.4

Still work to do...

# HMC with millions of parameters



J. Jasche

