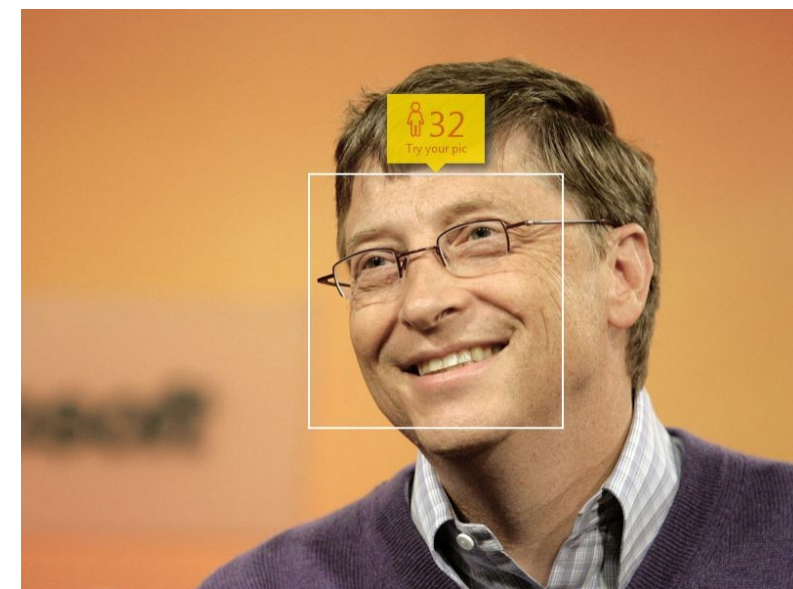# A Case Study of Machine Learning Hardware: Real-Time Source Separation using Markov Random Fields via Sampling-based Inference

## Glenn G. Ko and Rob A. Rutenbar
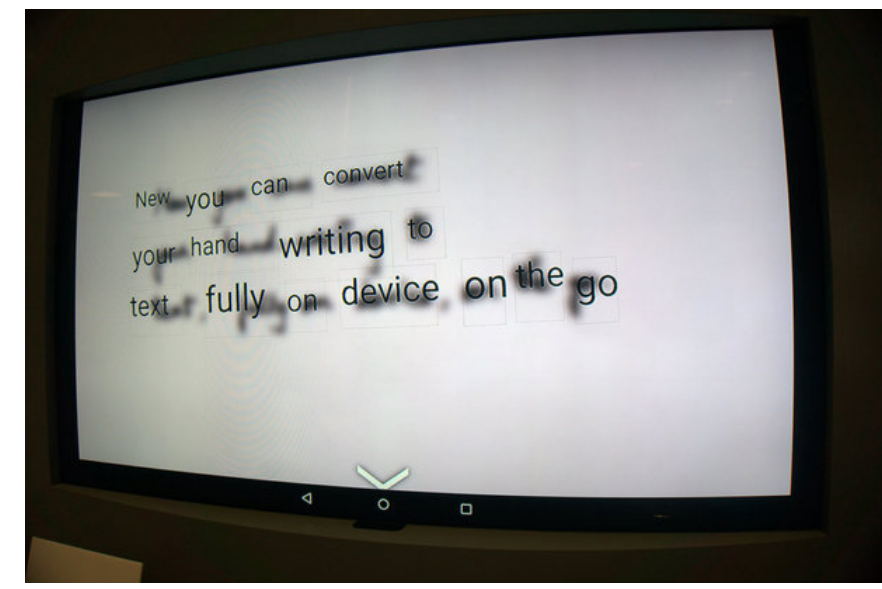## University of Illinois at Urbana-Champaign

## Motivation

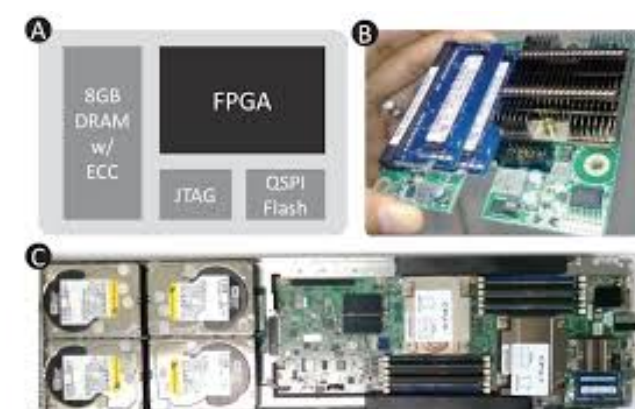### Why Machine Learning in Hardware?

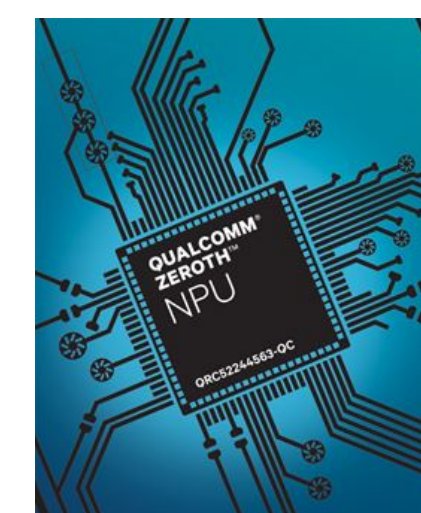Enterprise App: facial recognition

Mobile App: Handwriting Recognition

**Apps**

**Problems**

Big, complex, slow. Slow to run, slower (days/weeks) to train.

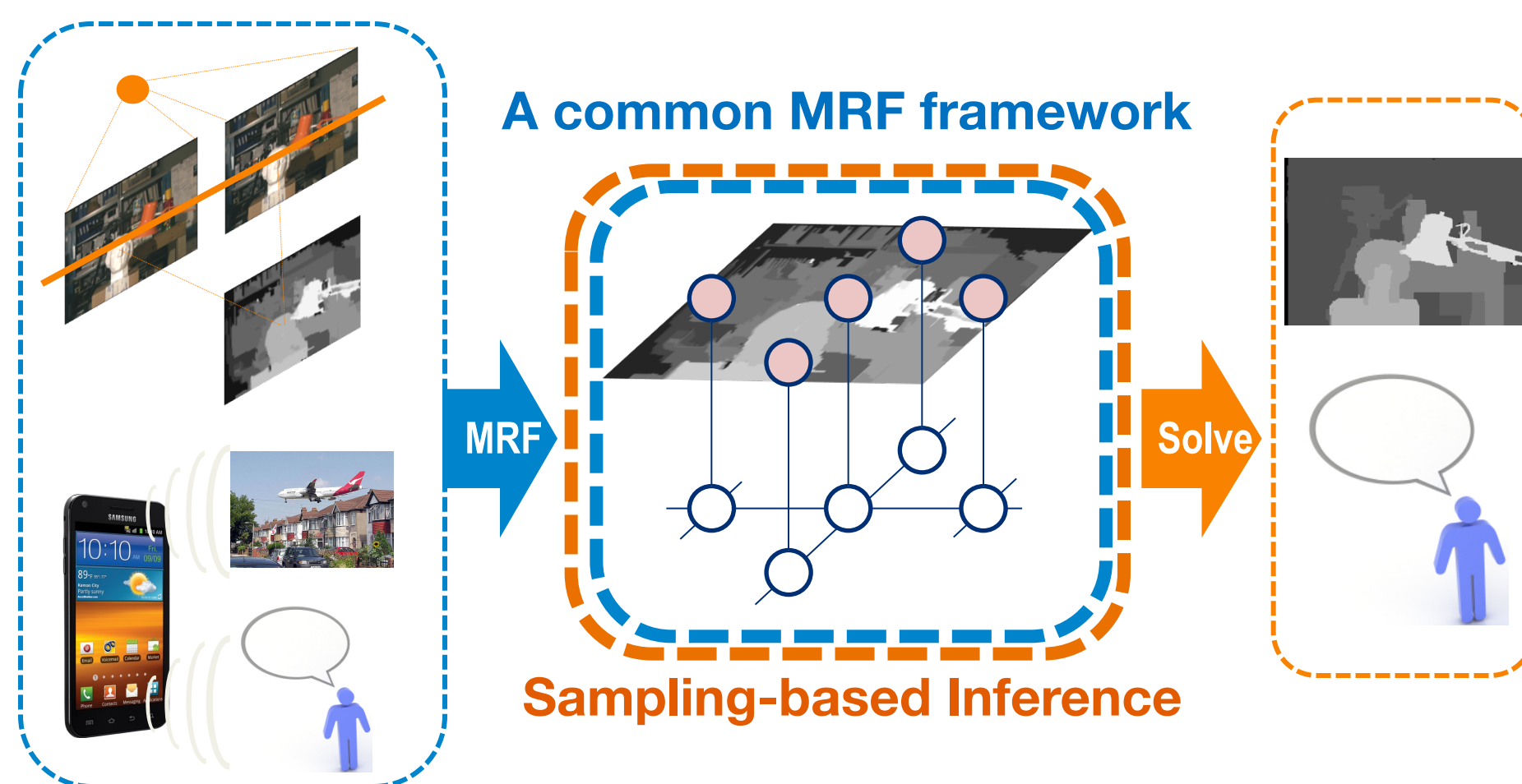Challenges involve running in real-time and low power
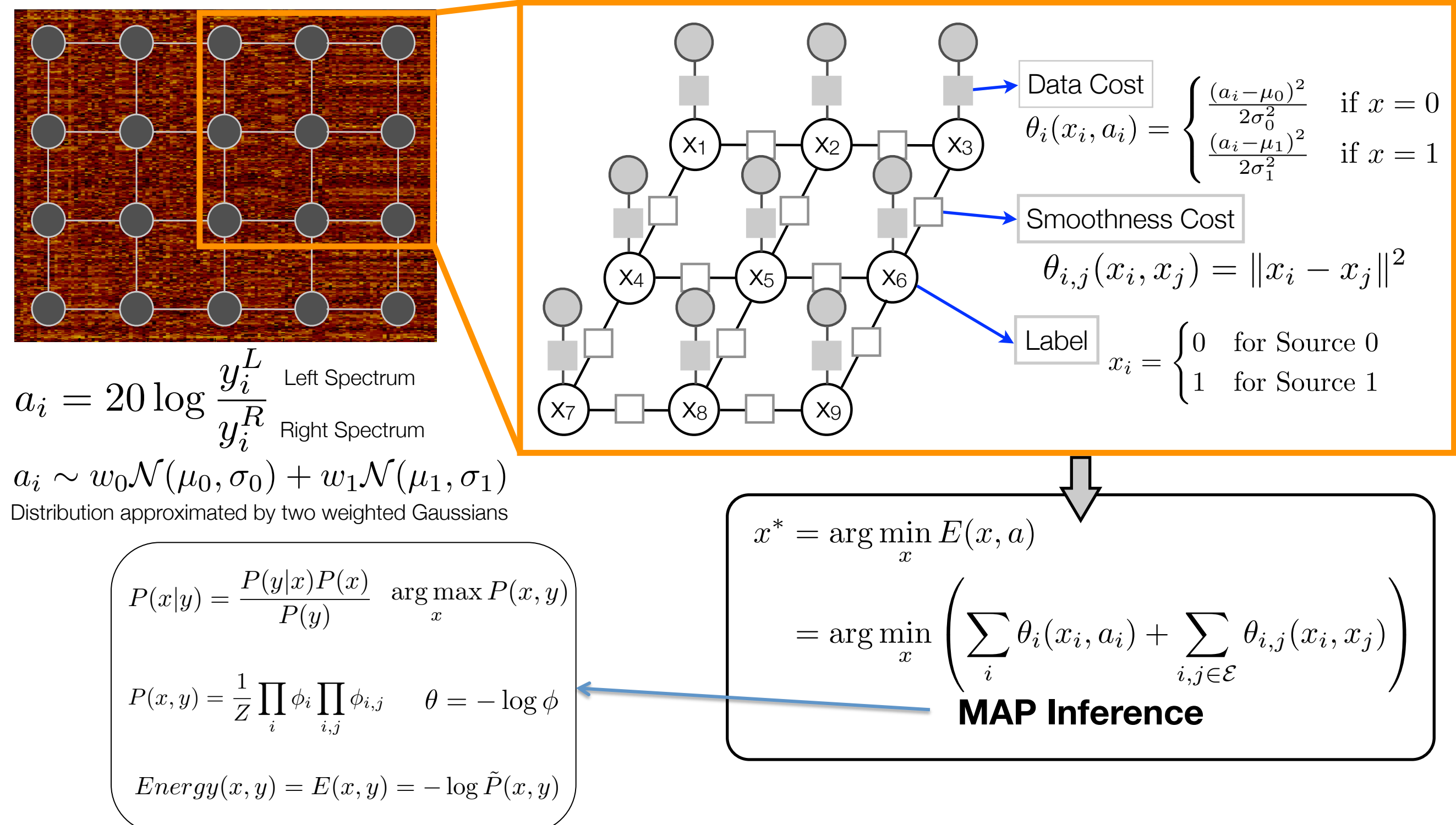
**Solutions**

Microsoft Catapult

Qualcomm Zeroth

## Objective

### Develop High Performance and Low-power Architectures for Inference on Probabilistic Graphical Models

- MRF is a general framework to solve ML inference applications
- An accelerator can possibly boost up all apps in this category

A common MRF framework

MRF → Solve

Sampling-based Inference

## Task: How to Separate Sounds?

Left

Right

Each point represents the energy of a (time, frequency) point

Captured Sound Mixtures

MRF Modeling

**Inference**

Binary Mask: which point belongs to which source

Undesired

Desired

Undesired: Source 0
Desired: Source 1

Separated Sources

## Mapping to Markov Random Field

$$\theta_i(x_i, a_i) = \begin{cases} \frac{(a_i - \mu_0)^2}{2\sigma_0^2} & \text{if } x = 0 \\ \frac{(a_i - \mu_1)^2}{2\sigma_1^2} & \text{if } x = 1 \end{cases}$$

Data Cost

Smoothness Cost $\theta_{i,j}(x_i, x_j) = \|x_i - x_j\|^2$

Label $x_i = \begin{cases} 0 & \text{for Source 0} \\ 1 & \text{for Source 1} \end{cases}$

$a_i = 20 \log \frac{y_i^L}{y_i^R}$   Left Spectrum / Right Spectrum

$a_i \sim w_0 \mathcal{N}(\mu_0, \sigma_0) + w_1 \mathcal{N}(\mu_1, \sigma_1)$
Distribution approximated by two weighted Gaussians

$$P(x|y) = \frac{P(y|x)P(x)}{P(y)} \quad \arg\max_x P(x, y)$$

$$P(x,y) = \frac{1}{Z} \prod_i \phi_i \prod_{i,j} \phi_{i,j} \quad \theta = -\log\phi$$

$$Energy(x,y) = E(x,y) = -\log \tilde{P}(x,y)$$

$$x^* = \arg\min_x E(x, a)$$

$$= \arg\min_x \left( \sum_i \theta_i(x_i, a_i) + \sum_{i,j \in \mathcal{E}} \theta_{i,j}(x_i, x_j) \right)$$

**MAP Inference**

## Inference using Gibbs Sampler

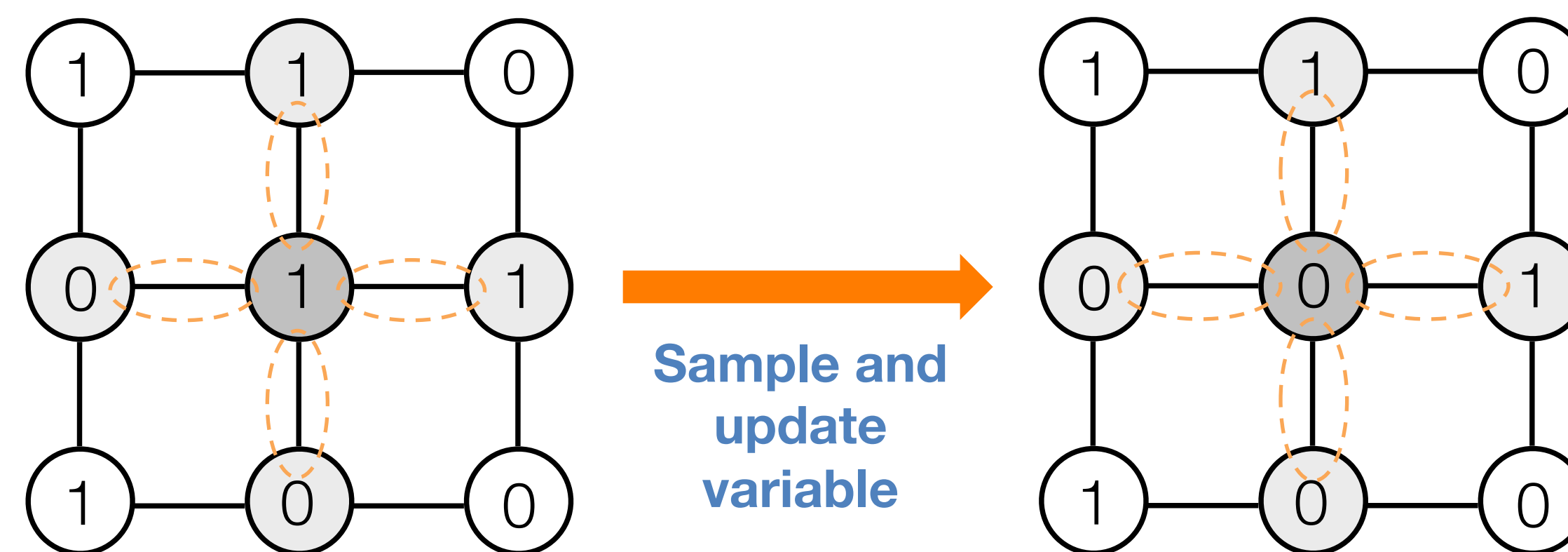**Algorithm 1** Gibbs sampling algorithm
1: **procedure** GIBBSSAMPLING(x)
2:   **for** t = 1 to T **do**
3:    **for** i = 1 to n **do**
4:     $x_i^{(t+1)} \sim P(x_i | x_1^{(t+1)}, ..., x_{i-1}^{t+1}, x_{i+1}^t, ..., x_n^{(t)})$
5:    **end for**
6:   **end for**
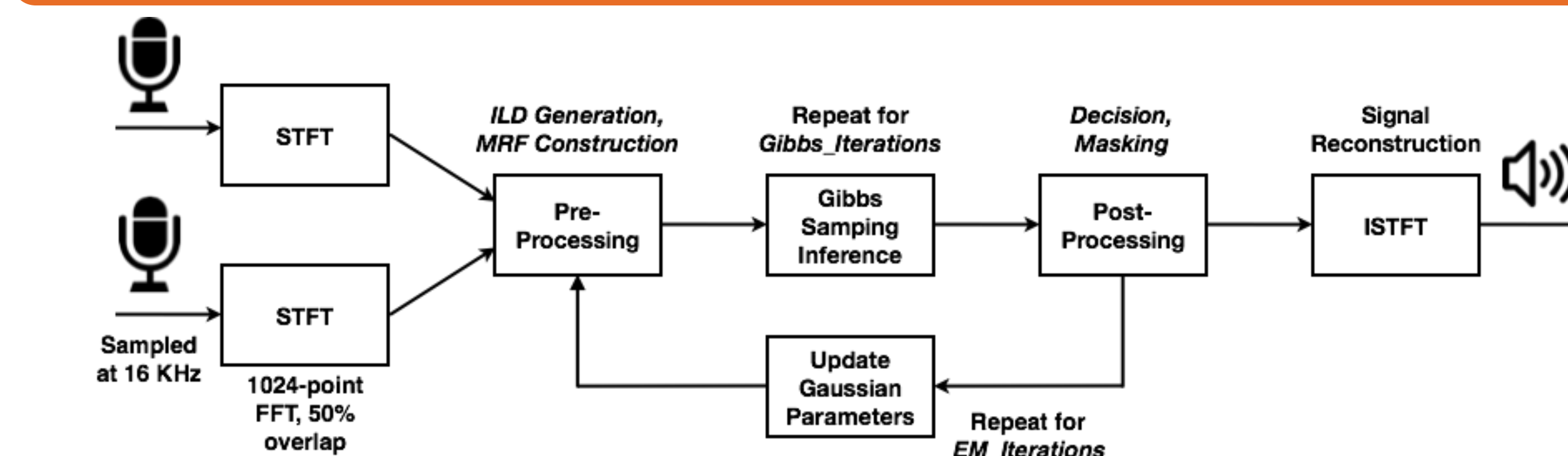7:   **return** samples
8: **end procedure**

Conditional Independence via Local Markov Property

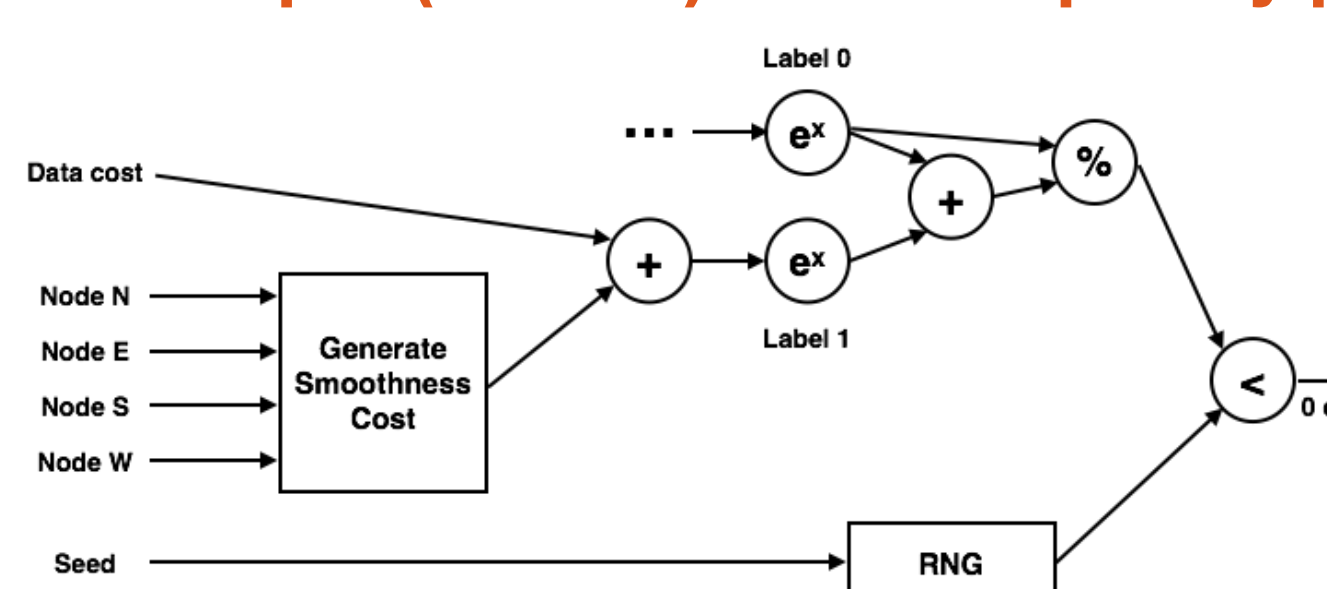$\forall x \in V : x \perp V \setminus x \mid neighbors(x)$

**Local computation: Only 4 data access for adjacent nodes**

Sample and update variable

## MCMC-EM Parameter Estimation

STFT

STFT

Sampled at 16 KHz

1024-point FFT, 50% overlap

ILD Generation, MRF Construction

Pre-Processing

Repeat for Gibbs_Iterations

Gibbs Samping Inference

Decision, Masking

Post-Processing

Signal Reconstruction

ISTFT

Update Gaussian Parameters

Repeat for EM_Iterations

**4 iterations of EM estimation and Gibbs sampling inference for MRF of 8 time pts (256 ms) x 513 frequency pts**
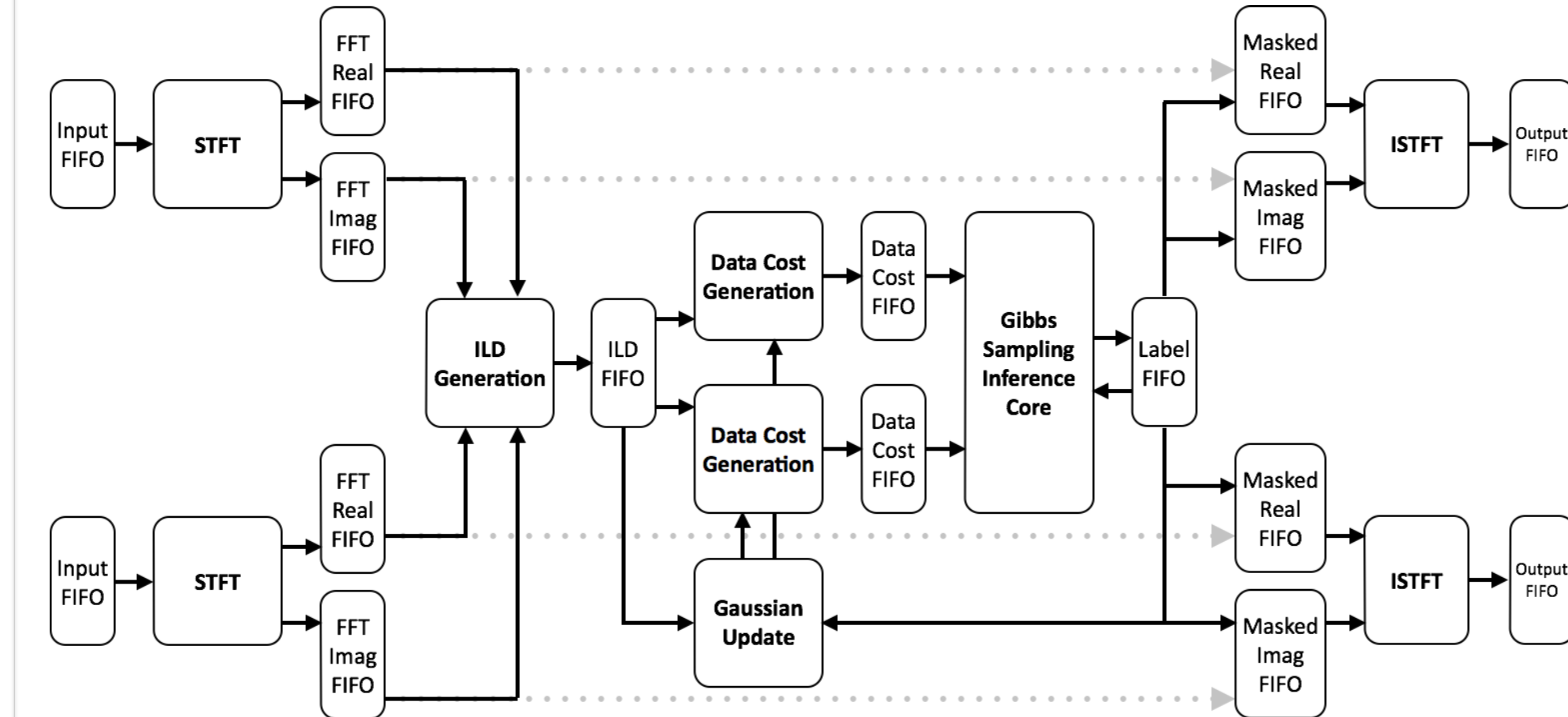
**Algorithm 1** Creating masks for source separation via MCMC-EM
1: **procedure** SOURCESEPARATIONMASK($A_i, x$)
2:   **for** each EM iteration **do**
3:    ConstructMRF
4:    (E-Step) GibbsSampling:
5:    **for** t = 1 to max iteration T **do**
6:     **for** each node i = 1 to I **do**
7:      $x_i^{(t+1)} \sim P(x_i | x_1^{(t+1)}, ..., x_{i-1}^{t+1}, x_{i+1}^t, ..., x_n^{(t)})$
8:     **end for**
9:    **end for**
10:    (M-Step) UpdateGaussians $\mu_{0,1}$ and $\sigma_{0,1}$
11:   **end for**
12:   **return** x
13: **end procedure**

Data cost

Node N
Node E
Node S
Node W

Generate Smoothness Cost

Seed

RNG

Label 0

Label 1

0 or 1

## Scalable Iterative Architecture

Input FIFO → STFT → FFT Real FIFO / FFT Imag FIFO → ILD Generation → ILD FIFO → Data Cost Generation → Data Cost FIFO → Gibbs Sampling Inference Core → Label FIFO → Masked Real FIFO / Masked Imag FIFO → ISTFT → Output FIFO

Gaussian Update

- Gibbs sampling inference core is a parameterized pipeline that can be easily extended to multiple parallel pipelines
- Uses inference results from previous frame for faster convergence

## Hardware Results

**Latency Requirements**
- ITU requirement 200 ms – LTE latency 160ms = 40ms: need < 40ms

**Software References**
- Intel Xeon X6550: ~64ms latency = too slow!
- ARM Cortex-A9 takes 23.32s to run 4s of audio with est. peak power of 3.67 Watts = slower and too power much for mobile!
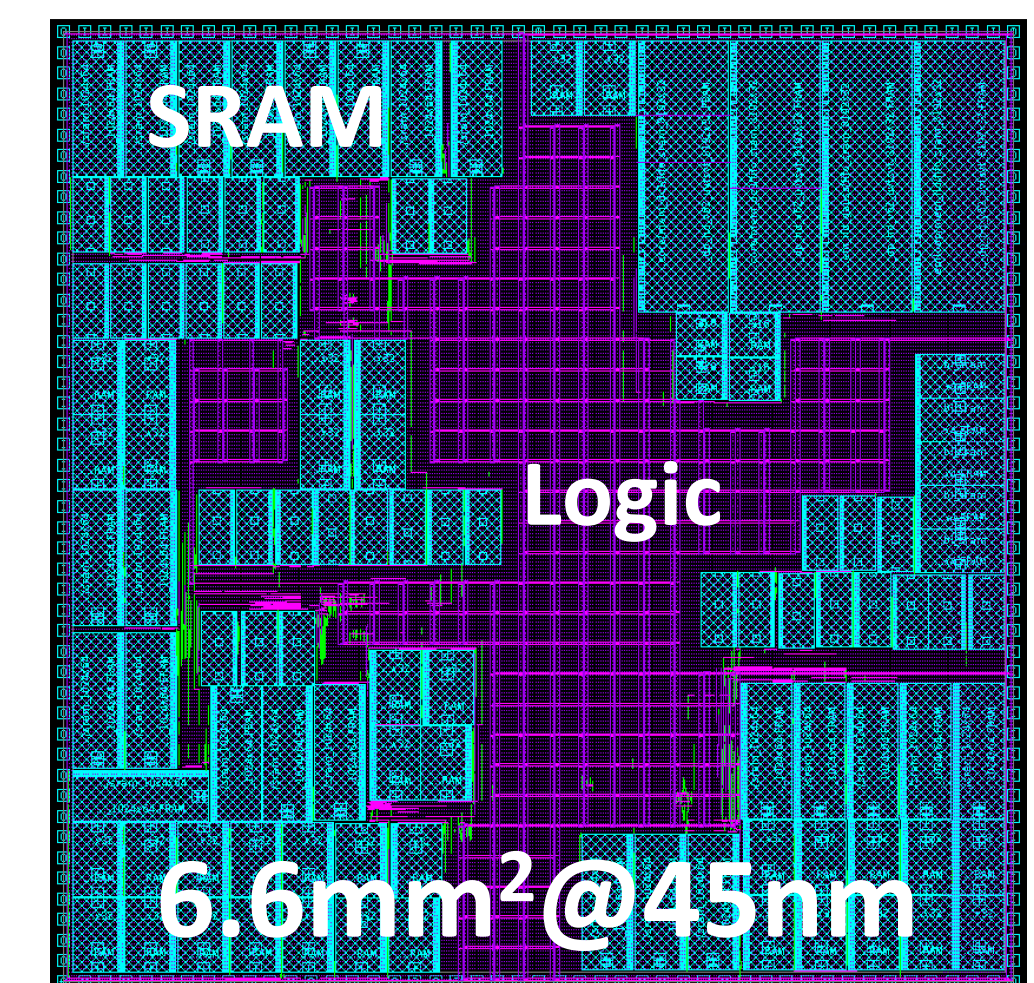
**FPGA Platform**
- Convey HC-1 w/ Xilinx Virtex5 FPGAs
- 150 MHz, 207 KB SRAM
- 6.7 dB Signal-to-Distortion-Ratio
- **1.6 ms real-time latency, 20X speedup**

**Virtual ASIC Design @ 45nm**
- Quite small < 10 M gates
- 469 mW at 150 MHz (404 mW from SRAM)
- 70 mW at 20 MHz (meets latency)
- **52X reduced power vs. ARM Cortex-A9**

| FPGA Resource | FPGA Utilization |
|---|---|
| Slice Register | 101314 / 207360 (48%) |
| Slice LUT | 90019 / 207360 (43%) |
| Slice LUT FF | 119418 / 207360 (57%) |
| BRAM | 115 / 288 (40%) |
| DSP | 36 / 192 (18%) |

SRAM

Logic

6.6mm²@45nm

## Summary

- A working HW implementation of source separation in a mobile form-factor using MRFs and Gibbs inference
- FPGA implementation running at **150 MHz** with **207 KB of RAM**, 64-bit memory width and 96 Kb/s of bandwidth requirement which is mobile-friendly
- Good **Signal-to-Distortion Ratio of 6.7 dB** with better auditory performance and very low **1.601ms latency**, a **20X speedup** from input to output sound stream, well below latency req for human hearing on mobile phone
- Virtual ASIC design (45 nm) **power estimate 70 mW running at 20 MHz to meet latency req, a 52X power reduction** vs. ARM Cortex-A9 mobile software reference

STARnet

SONIC