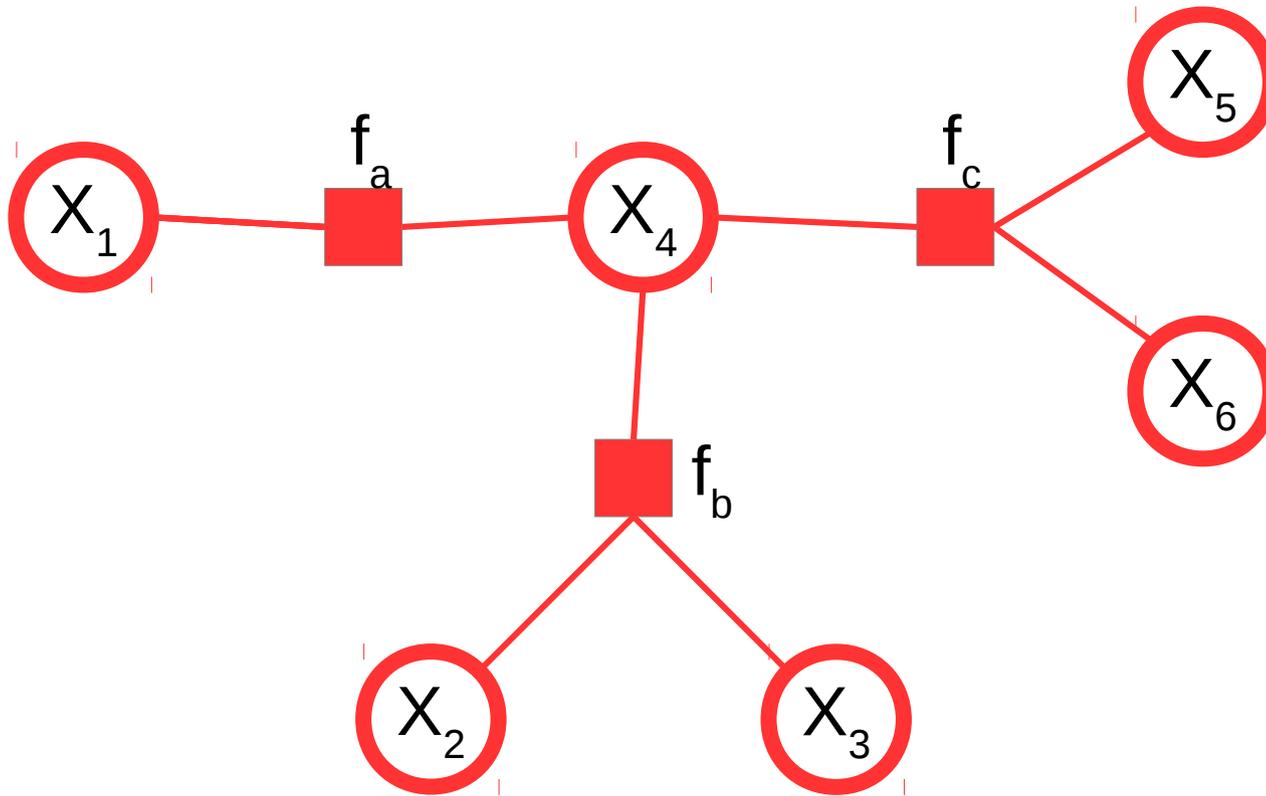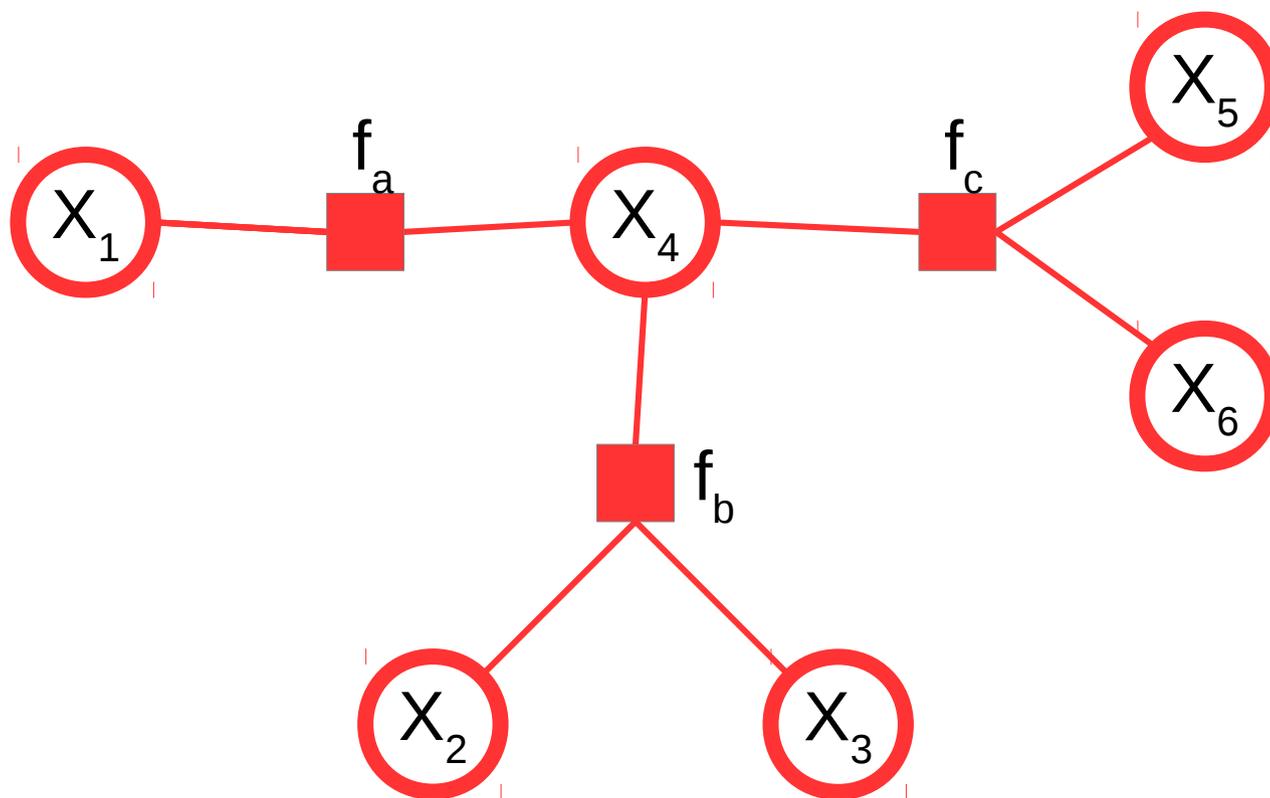# Inference in factor trees
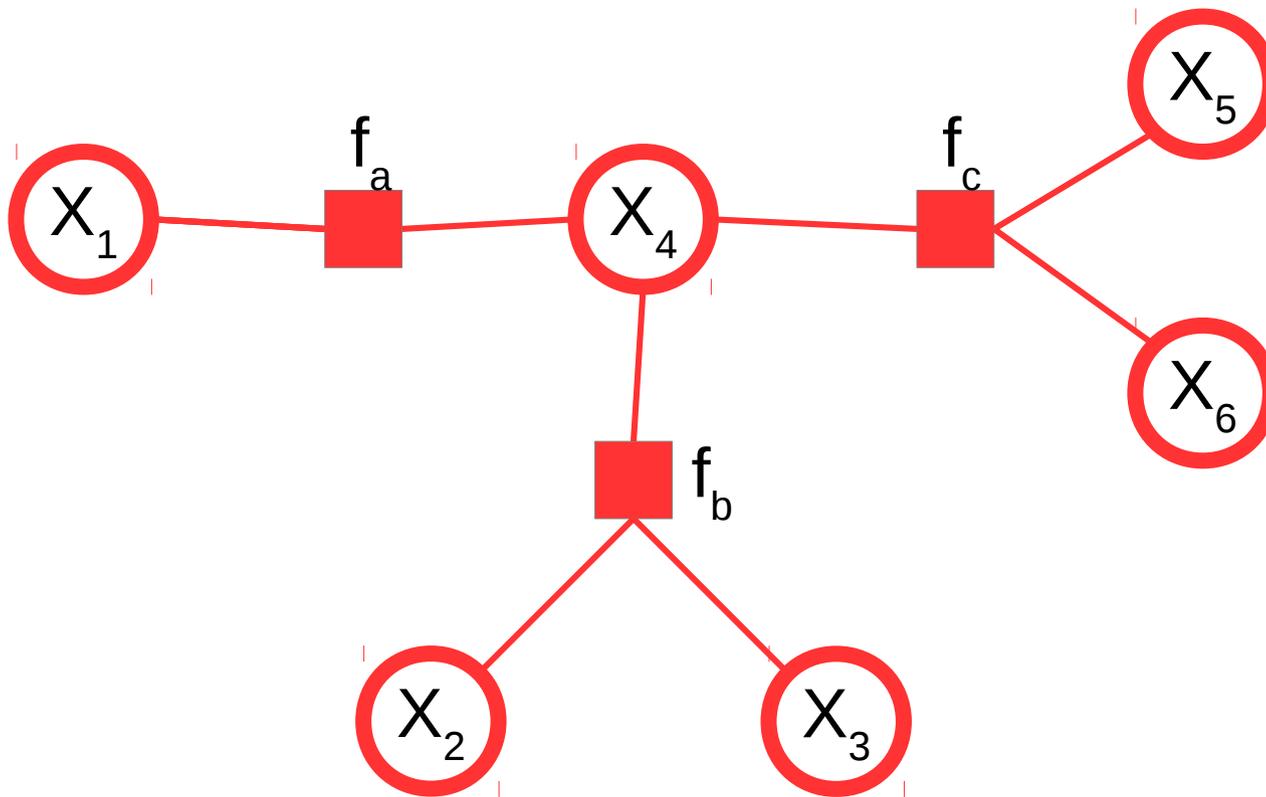
**Carlos Cotrini**
**November 3, 2017**

**Probabilistic foundations of artificial intelligence**

$$P(X_1 = \hat{x_1}, \ldots, X_6 = \hat{x_6}) =$$

$$\frac{1}{Z} f_a(\hat{x_1}, \hat{x_4}) f_b(\hat{x_2}, \hat{x_3}, \hat{x_4}) f_c(\hat{x_4}, \hat{x_5}, \hat{x_6})$$

3

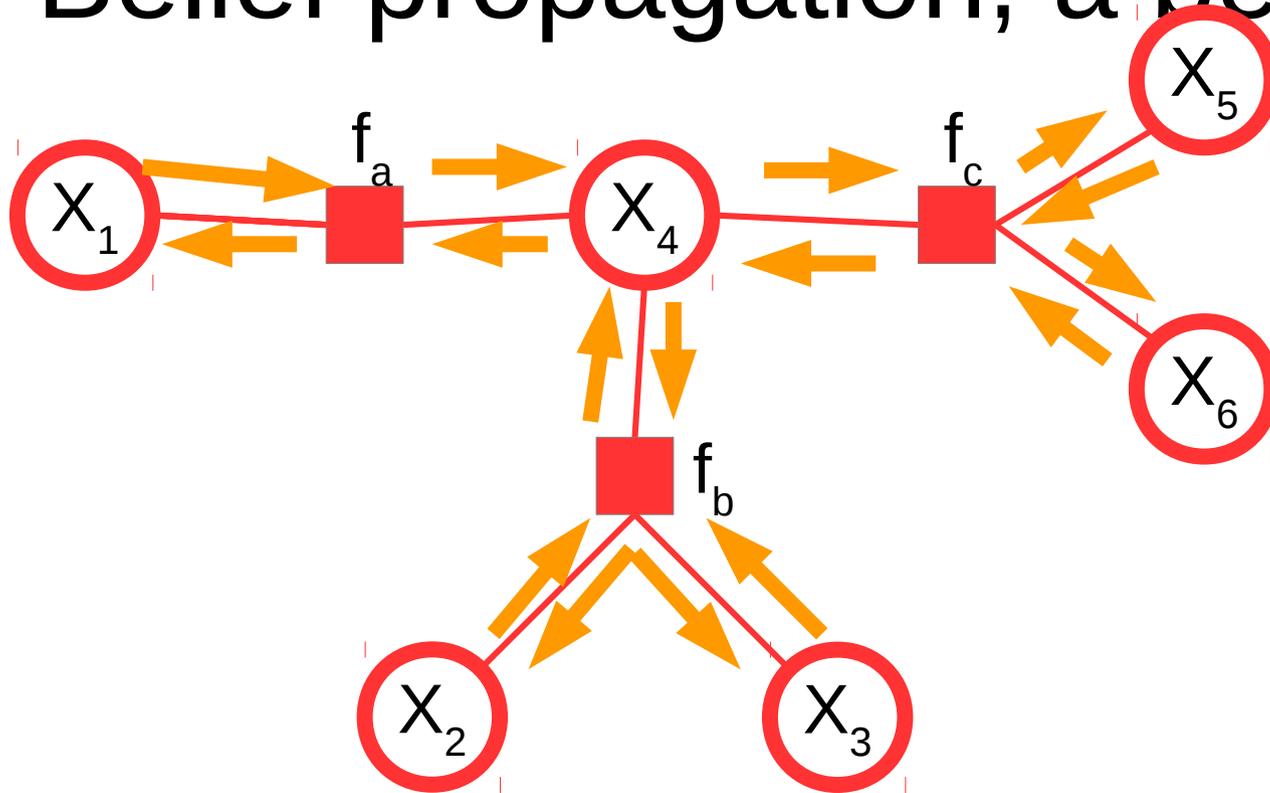$$P(X_1 = \hat{x}_1, \ldots, X_6 = \hat{x}_6) =$$

$$\frac{1}{Z} f_a(\hat{x}_1, \hat{x}_4) f_b(\hat{x}_2, \hat{x}_3, \hat{x}_4) f_c(\hat{x}_4, \hat{x}_5, \hat{x}_6)$$

How do we compute $P(X_5 = \hat{x}_5)$?

4

# Naive method

$$P(X_5 = \hat{x}_5) =$$
$$\sum_{\hat{x}_1, \hat{x}_2, \hat{x}_3, \hat{x}_4, \hat{x}_6} P(X_1 = \hat{x}_1, X_2 = \hat{x}_2, X_3 = \hat{x}_3, X_4 = \hat{x}_4, X_6 = \hat{x}_6).$$

# Belief propagation, a better method.

# Belief propagation, a better method.

$$\mu^{(t)}_{X \to f}(\hat{x}) = \prod_{f' \in N(X) \setminus \{f\}} \mu^{(t-1)}_{f' \to X}(\hat{x})$$

$$\mu^{(t)}_{f \to X}(\hat{x}) = \sum_{\mathbf{x}} f(\hat{\mathbf{x}}, \hat{x}) \prod_{X' \in N(f) \setminus \{X\}} \mu^{(t-1)}_{X' \to f}(\hat{x'})$$

- $X$: a node (i.e., a random variable).
- $f$: a factor.
- $\hat{x}$: a value in the range of $X$.
- $N(X)$: $X$'s neighbors.
- $N(f)$: $f$'s neighbors.
- $\hat{\mathbf{x}}$: a sequence of values in $f$'s domain.
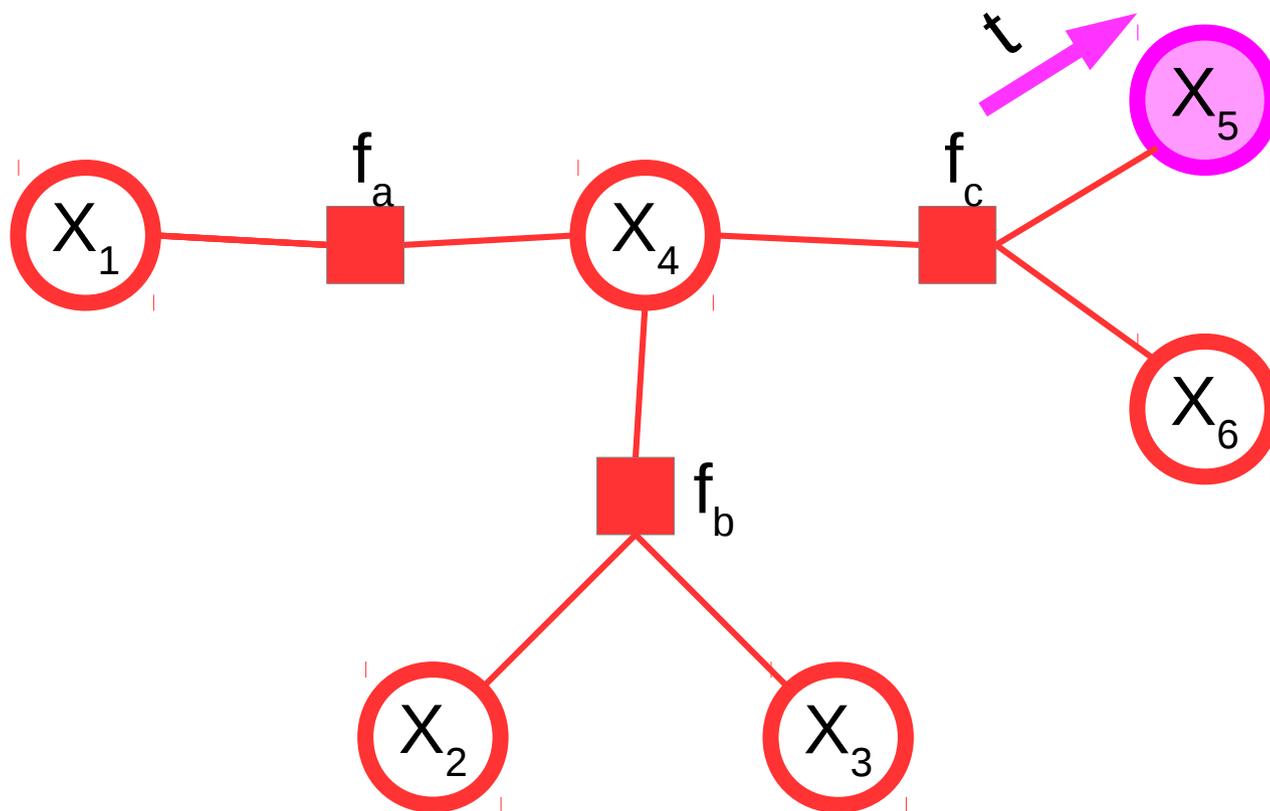
Initially,
$$\mu^{(0)}_{X \to f}(\hat{x}) = 1 \text{ and }$$
$$\mu^{(0)}_{f \to X}(\hat{x}) = 1.$$

Before we compute $P(X_5 = x_5)$, let's observe three useful insights about belief propagation in trees.

# First insight
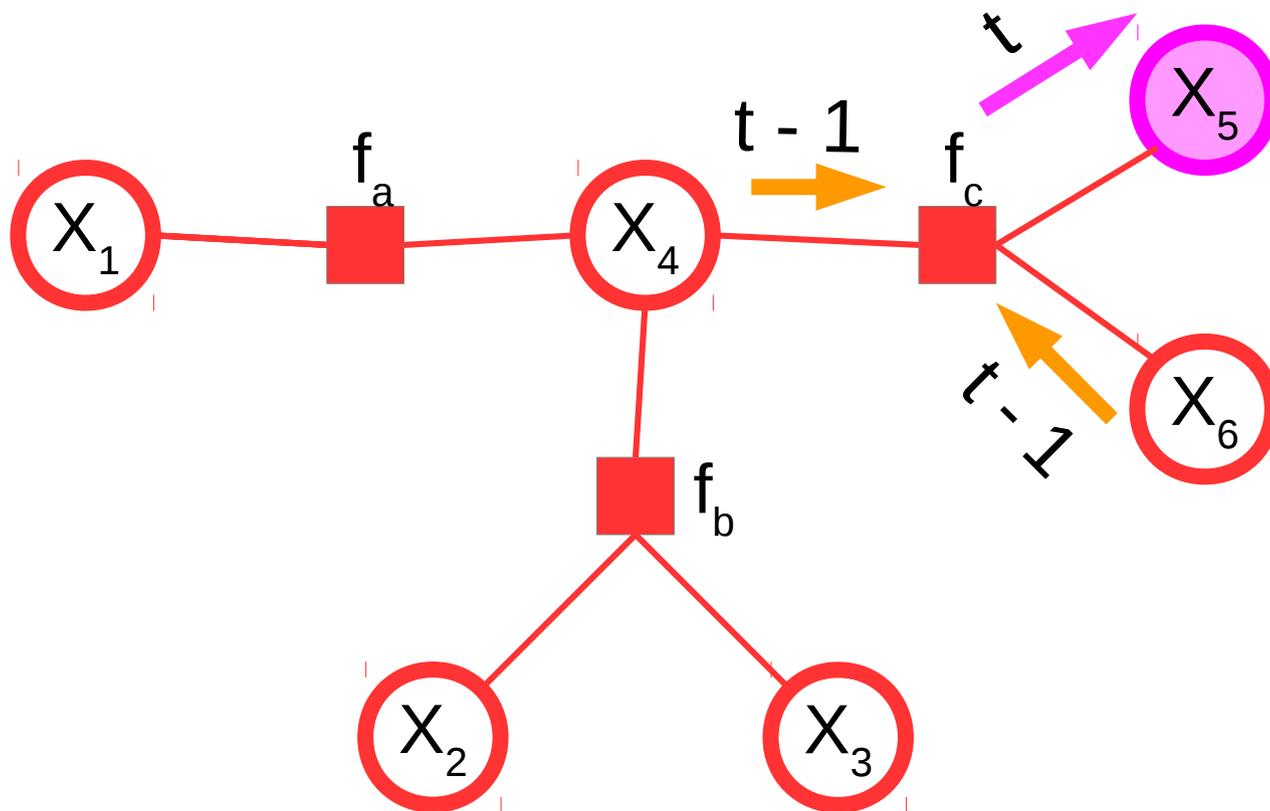
The messages needed to compute another message form a tree*
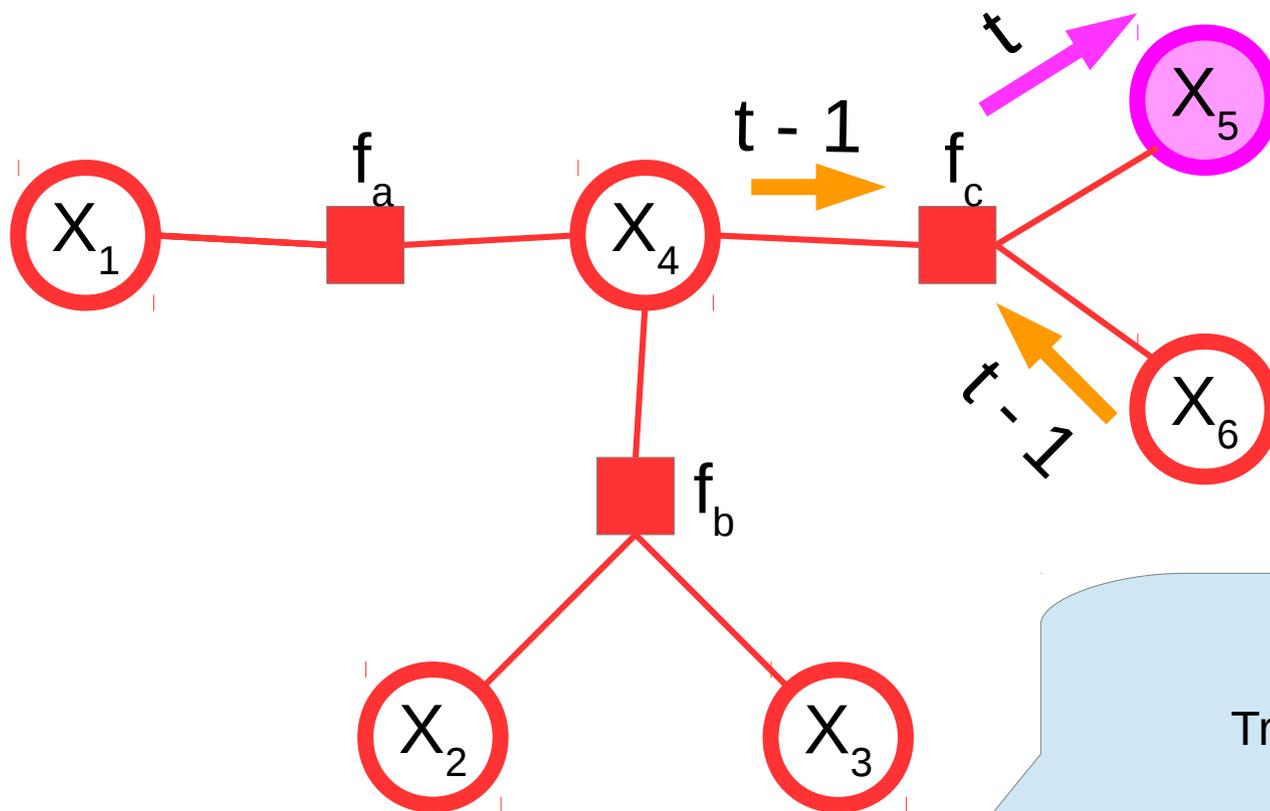
\* This only holds for factor trees!

$$\mu_{X \to f}^{(t)}(\hat{x}) = \prod_{f' \in N(X) \backslash \{f\}} \mu_{f' \to X}^{(t-1)}(\hat{x})$$

$$\mu_{f \to X}^{(t)}(\hat{x}) = \sum_{\mathbf{x}} f(\hat{\hat{\mathbf{x}}}, \hat{x}) \prod_{X' \in N(f) \backslash \{X\}} \mu_{X' \to f}^{(t-1)}(\hat{x}')$$

$$\mu_{X \to f}^{(t)}(\hat{x}) = \prod_{f' \in N(X) \setminus \{f\}} \mu_{f' \to X}^{(t-1)}(\hat{x})$$

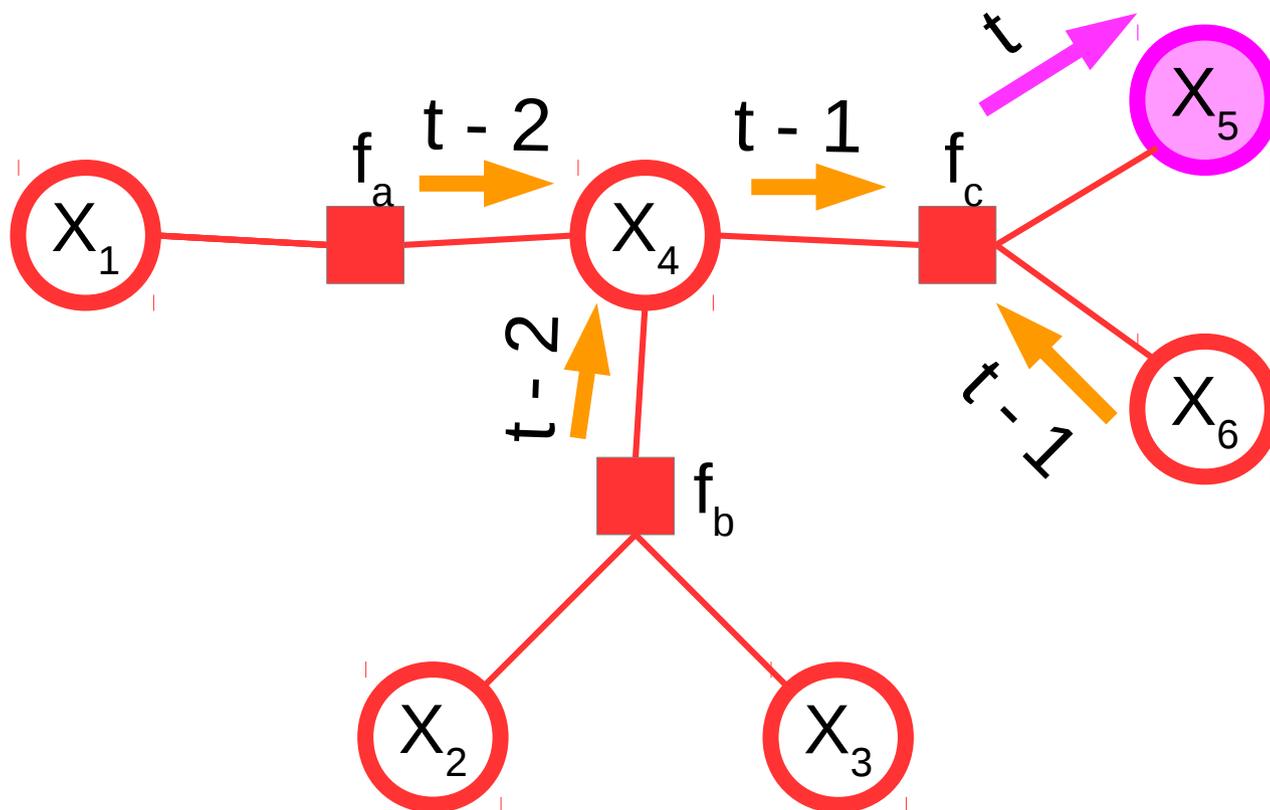$$\mu_{f \to X}^{(t)}(\hat{x}) = \sum_{\mathbf{x}} f(\hat{\hat{\mathbf{x}}}, \hat{x}) \prod_{X' \in N(f) \setminus \{X\}} \mu_{X' \to f}^{(t-1)}(\hat{x}')$$
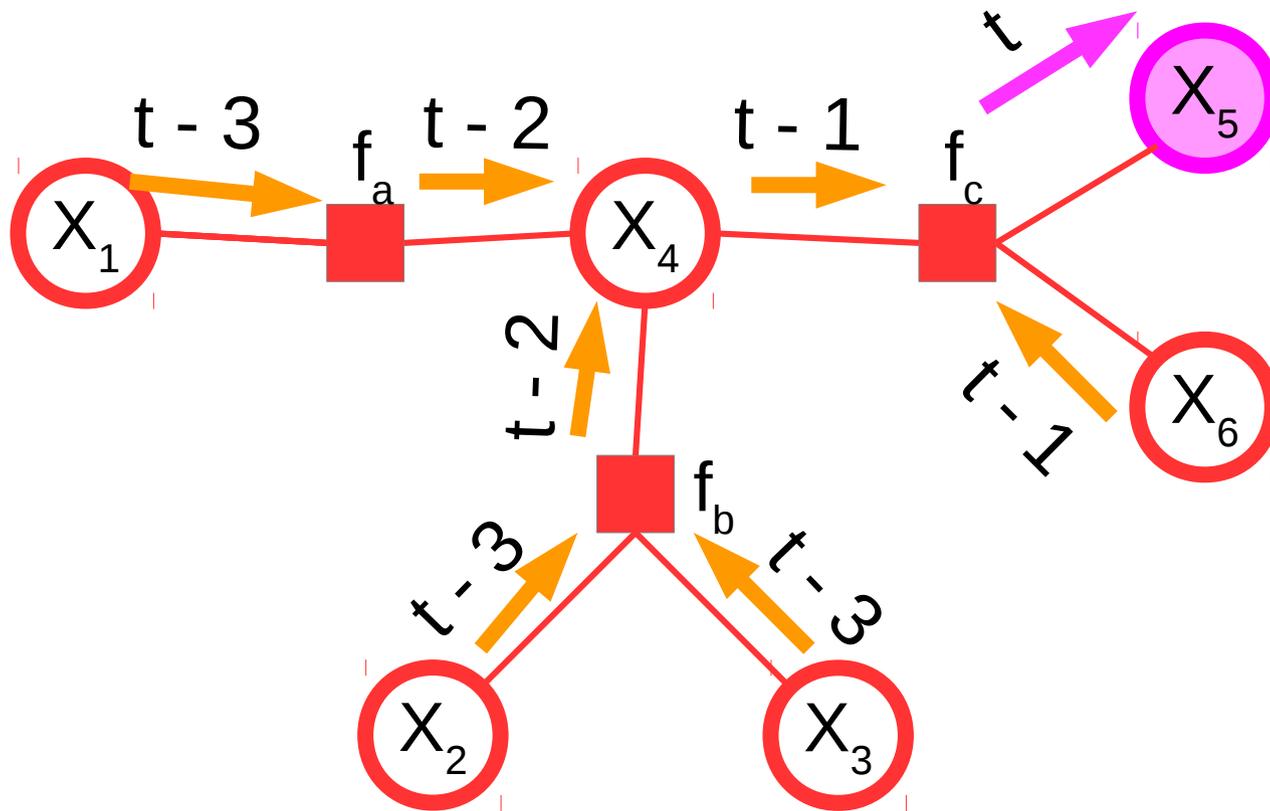
$$\mu_{X \to f}^{(t)}(\hat{x}) = \prod_{f' \in N(X) \setminus \{f\}} \mu_{f' \to X}^{(t-1)}(\hat{x})$$

$$\mu_{f \to X}^{(t)}(\hat{x}) = \sum_{\mathbf{x}} f(\hat{\mathbf{x}}, \hat{x}) \prod_{X' \in N(f) \setminus \{X\}} \mu_{X' \to f}^{(t-1)}(\hat{x'})$$

$$\mu_{X \to f}^{(t)}(\hat{x}) = \prod_{f' \in N(X) \setminus \{f\}} \mu_{f' \to X}^{(t-1)}(\hat{x})$$
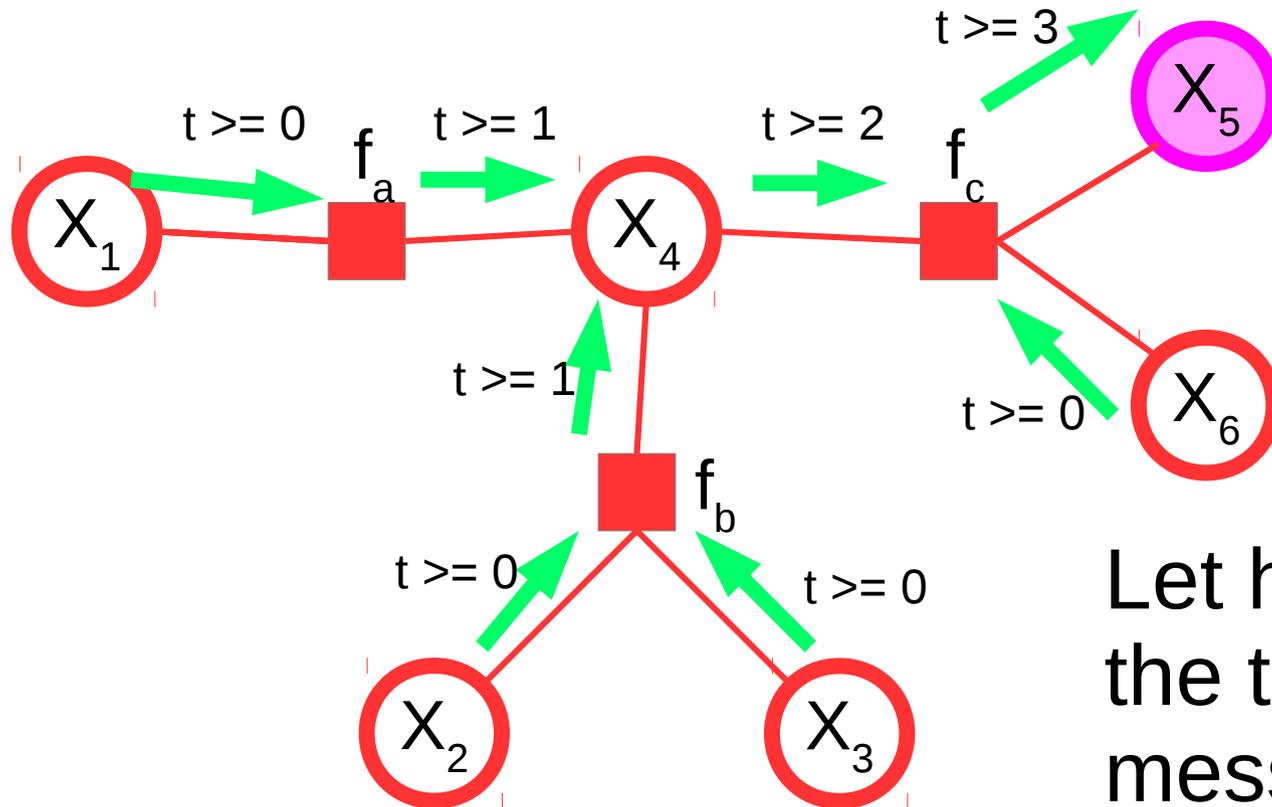
$$\mu_{f \to X}^{(t)}(\hat{x}) = \sum_{\mathbf{x}} f(\hat{\hat{\mathbf{x}}}, \hat{x}) \prod_{X' \in N(f) \setminus \{X\}} \mu_{X' \to f}^{(t-1)}(\hat{x'})$$

$$\mu_{X \to f}^{(t)}(\hat{x}) = \prod_{f' \in N(X) \setminus \{f\}} \mu_{f' \to X}^{(t-1)}(\hat{x})$$

$$\mu_{f \to X}^{(t)}(\hat{x}) = \sum_{\mathbf{x}} f(\hat{\hat{\mathbf{x}}}, \hat{x}) \prod_{X' \in N(f) \setminus \{X\}} \mu_{X' \to f}^{(t-1)}(\hat{x}')$$

# First insight

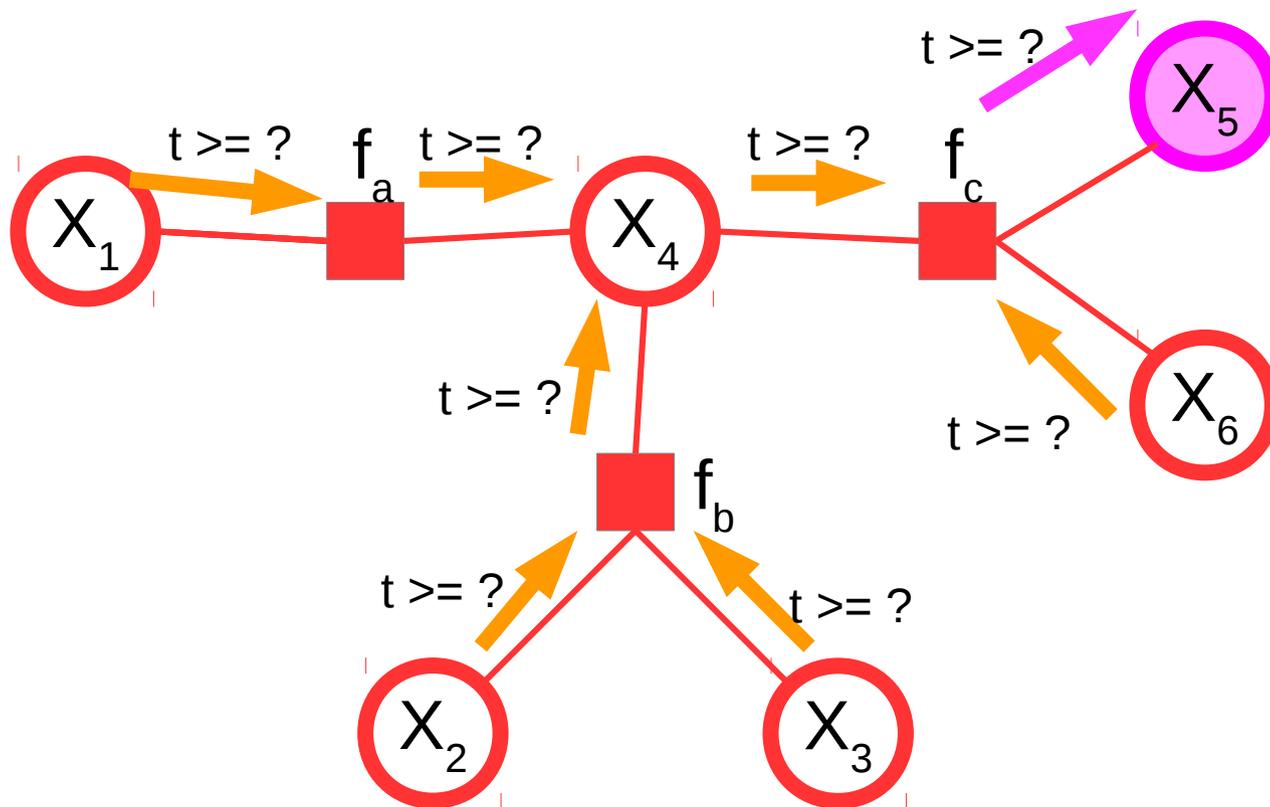The messages needed to compute another message form a tree*
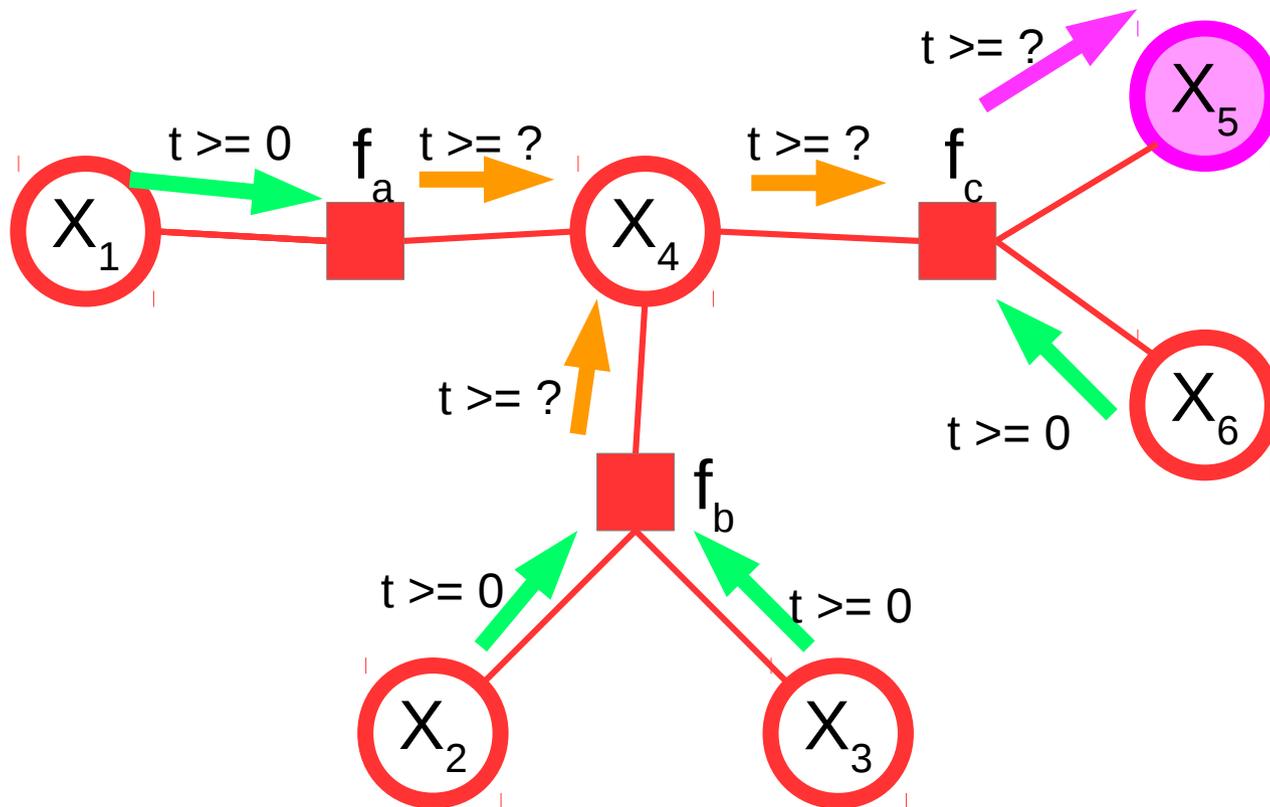
* This only holds for trees!

# Second insight



Let h be the height of the tree rooted at one message.

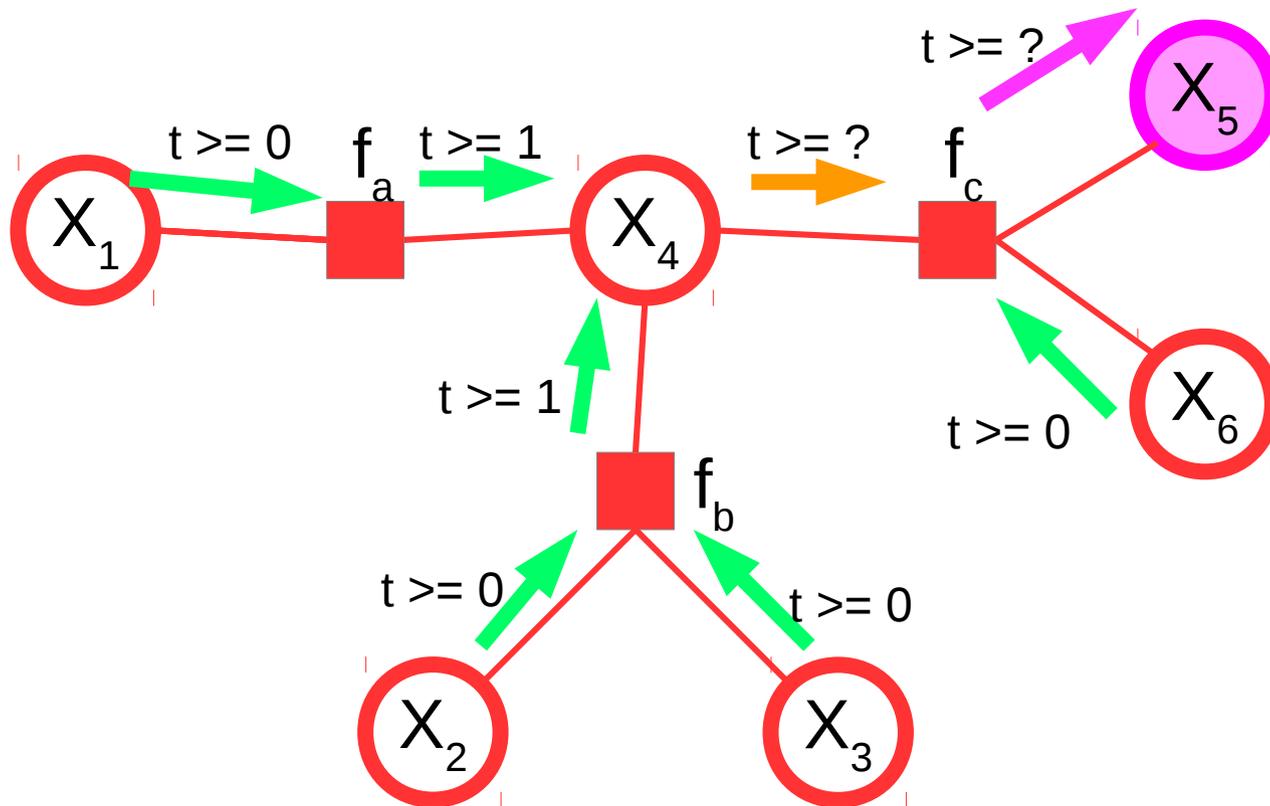- At iteration h of BP, the message reaches its final value.

$$\mu_{X \to f}^{(t)}(\hat{x}) = \prod_{f' \in N(X) \setminus \{f\}} \mu_{f' \to X}^{(t-1)}(\hat{x})$$

$$\mu_{f \to X}^{(t)}(\hat{x}) = \sum_{\mathbf{x}} f(\hat{\hat{\mathbf{x}}}, \hat{x}) \prod_{X' \in N(f) \setminus \{X\}} \mu_{X' \to f}^{(t-1)}(\hat{x'})$$
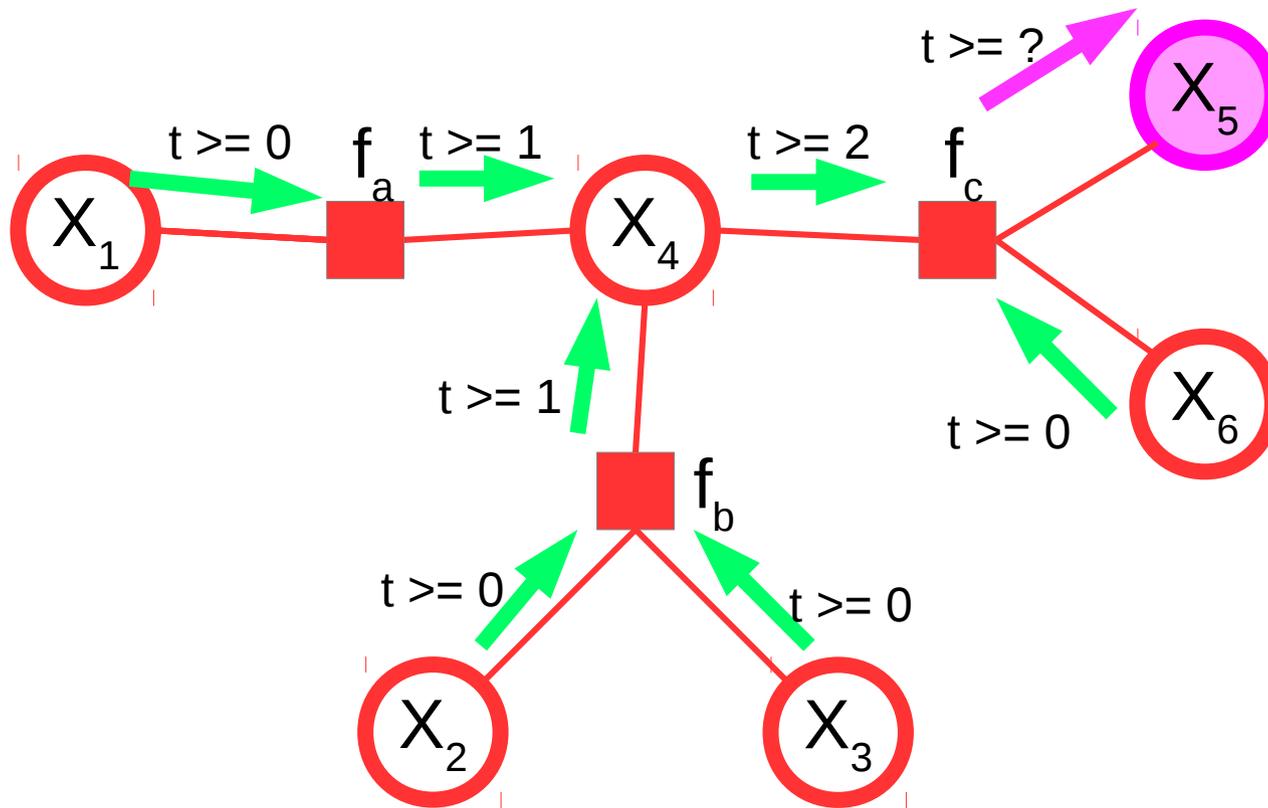
$$\mu_{X \to f}^{(t)}(\hat{x}) = \prod_{f' \in N(X) \setminus \{f\}} \mu_{f' \to X}^{(t-1)}(\hat{x})$$

$$\mu_{f \to X}^{(t)}(\hat{x}) = \sum_{\mathbf{x}} f(\hat{\hat{\mathbf{x}}}, \hat{x}) \prod_{X' \in N(f) \setminus \{X\}} \mu_{X' \to f}^{(t-1)}(\hat{x'})$$
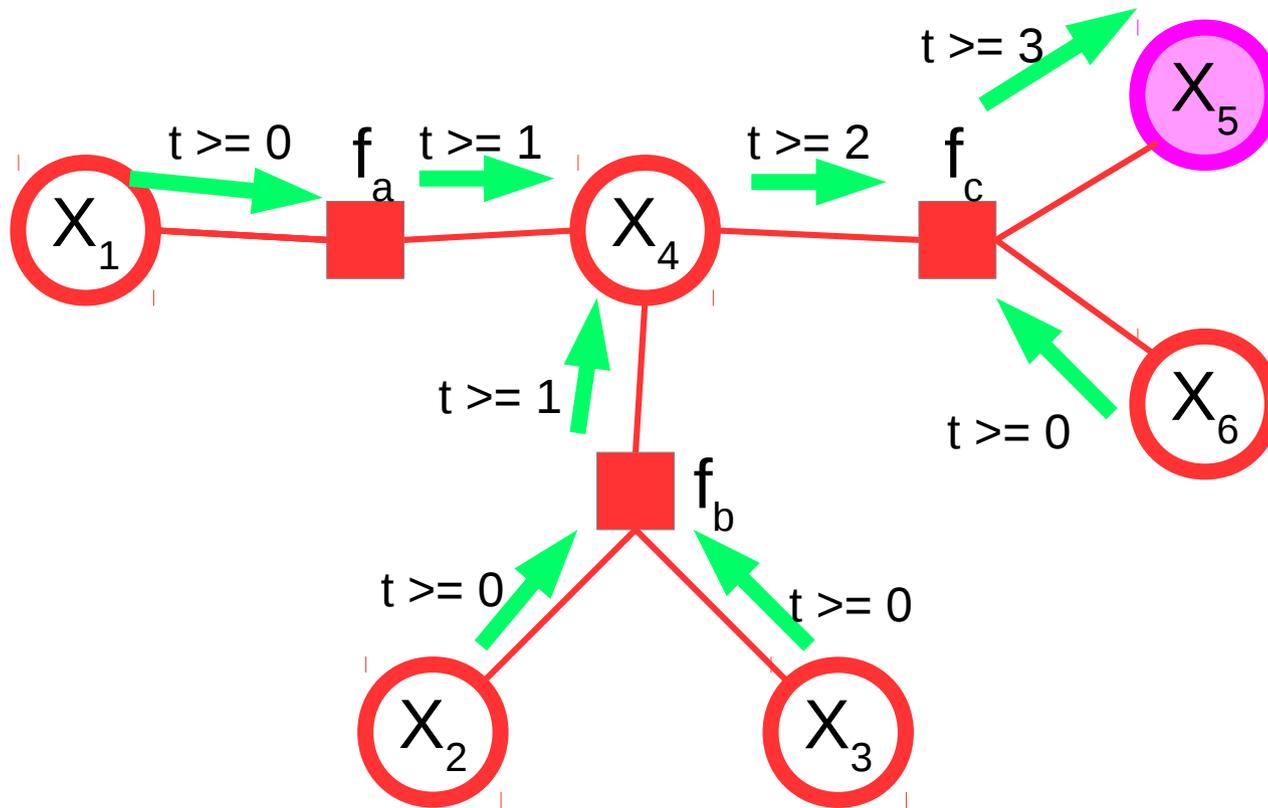
$$\mu_{X \to f}^{(t)}(\hat{x}) = \prod_{f' \in N(X) \setminus \{f\}} \mu_{f' \to X}^{(t-1)}(\hat{x})$$
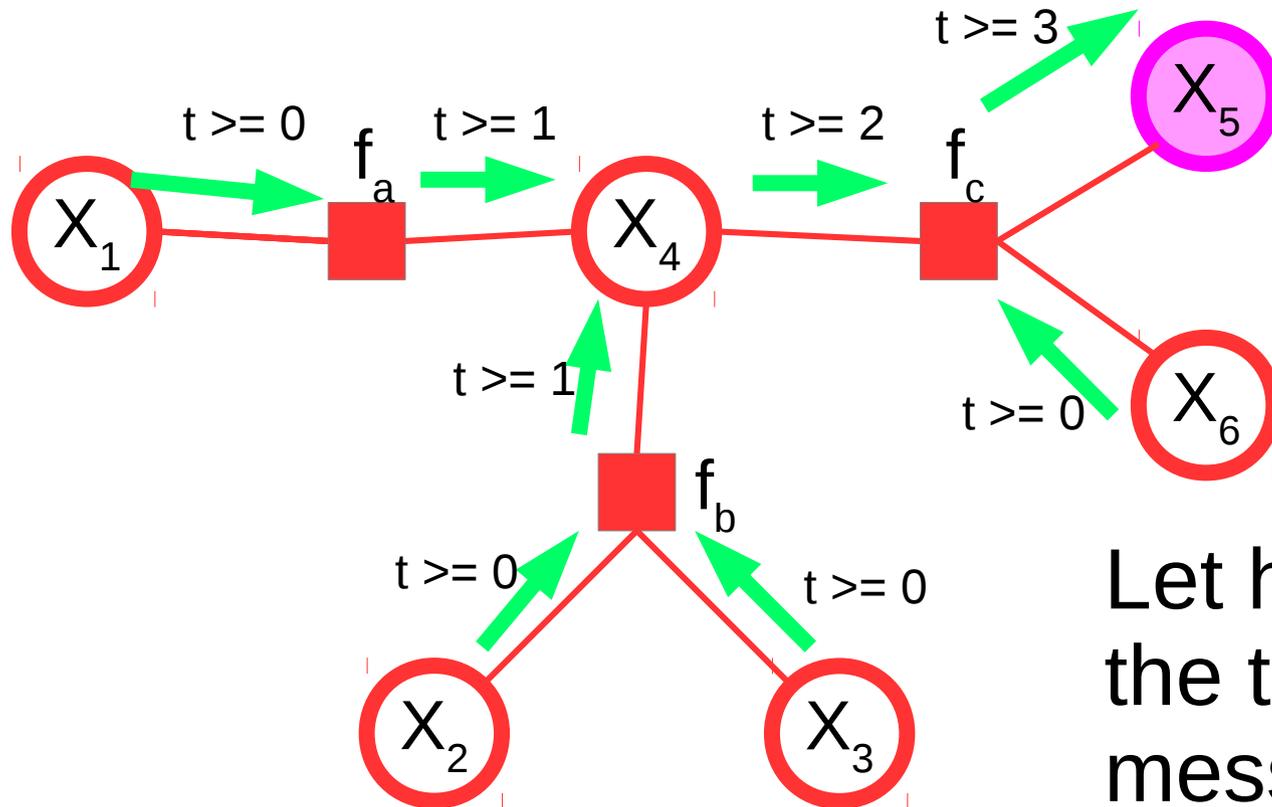
$$\mu_{f \to X}^{(t)}(\hat{x}) = \sum_{\mathbf{x}} f(\hat{\hat{\mathbf{x}}}, \hat{x}) \prod_{X' \in N(f) \setminus \{X\}} \mu_{X' \to f}^{(t-1)}(\hat{x'})$$

$$\mu_{X \to f}^{(t)}(\hat{x}) = \prod_{f' \in N(X) \setminus \{f\}} \mu_{f' \to X}^{(t-1)}(\hat{x})$$

$$\mu_{f \to X}^{(t)}(\hat{x}) = \sum_{\mathbf{x}} f(\hat{\hat{\mathbf{x}}}, \hat{x}) \prod_{X' \in N(f) \setminus \{X\}} \mu_{X' \to f}^{(t-1)}(\hat{x'})$$

$$\mu^{(t)}_{X \to f}(\hat{x}) = \prod_{f' \in N(X) \backslash \{f\}} \mu^{(t-1)}_{f' \to X}(\hat{x})$$

$$\mu^{(t)}_{f \to X}(\hat{x}) = \sum_{\mathbf{x}} f(\hat{\hat{\mathbf{x}}}, \hat{x}) \prod_{X' \in N(f) \backslash \{X\}} \mu^{(t-1)}_{X' \to f}(\hat{x'})$$
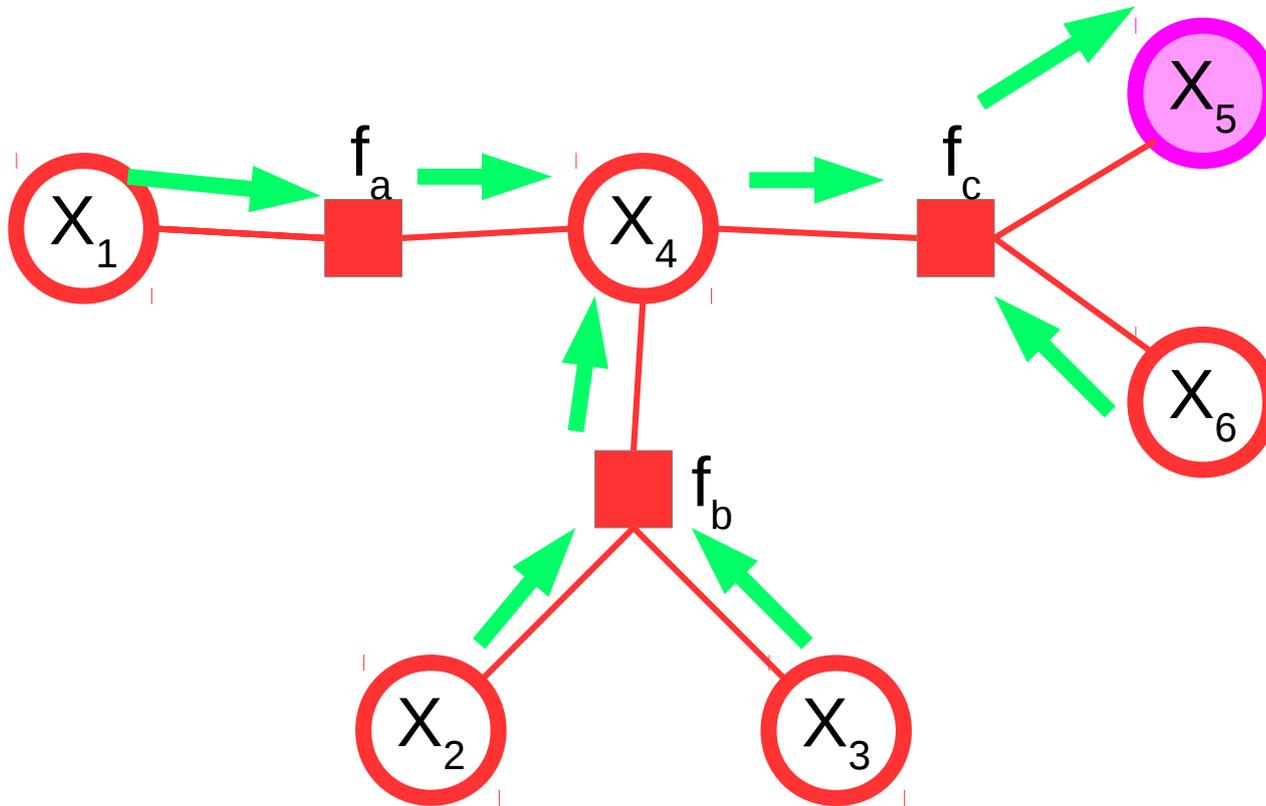
# Second insight

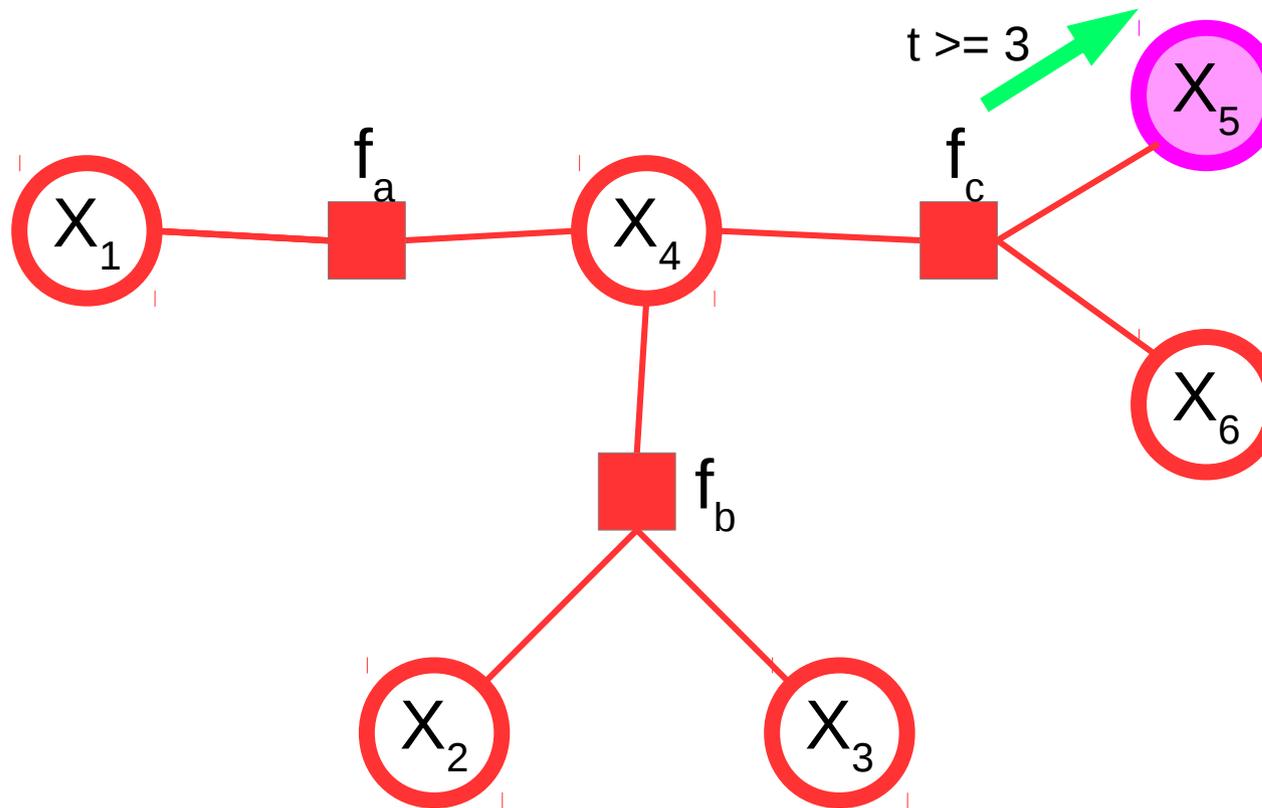

Let h be the height of the tree rooted at one message.

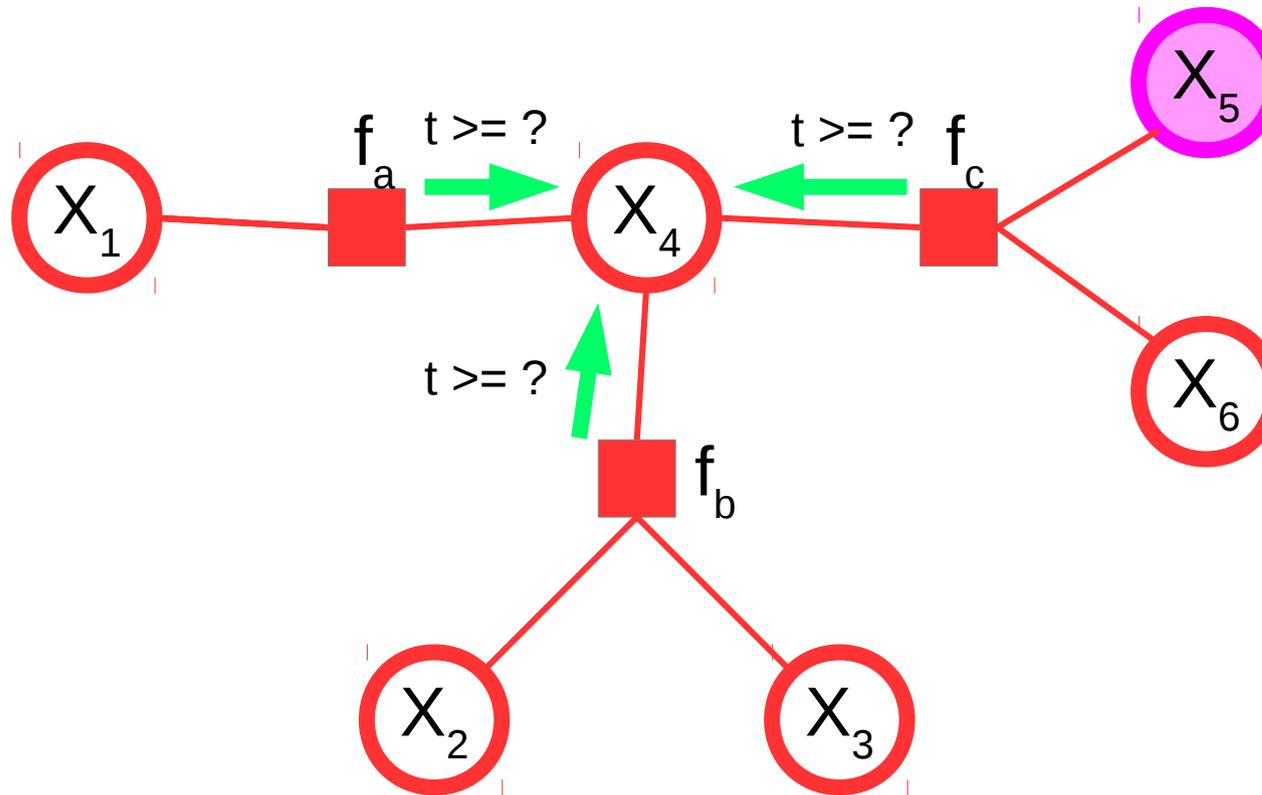- At iteration h, the message reaches its final value.

# Third insight



Let $t_0$ be the time at which all messages have reached their true value. For any $t \geq t_0$, $P(X = \hat{x}) = \frac{1}{Z} \prod_{f \in N(X)} \mu_{f \to X}^{(t)}(\hat{x})$.
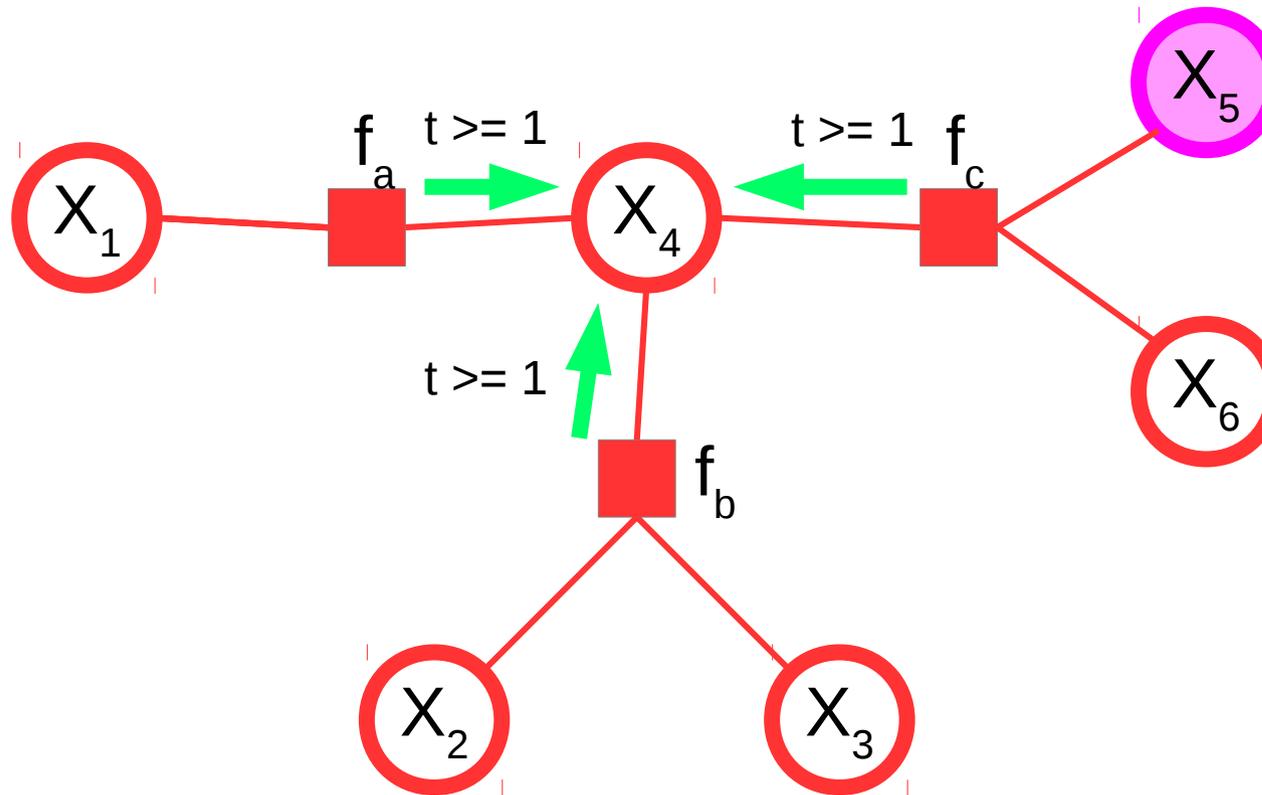
# $P(X_5 = x_5)$

# $P(X_4 = x_4)$

# $P(X_4 = x_4)$

# Rejection sampling and MCMC

**Carlos Cotrini**
**November 3, 2017**

**Probabilistic foundations of artificial intelligence**

# Computing expected loss from a burglary



# What is the expected total value of items stolen?

# Computing expected loss from a burglary

# Given…

- Col, a collection of objects.
- Cap, the capacity of the bag.

- Let $\mathsf{FIT} := \{A \subseteq \mathsf{Col} \mid \sum_{b \in A} weight(b) \leq \mathsf{Cap}\}$

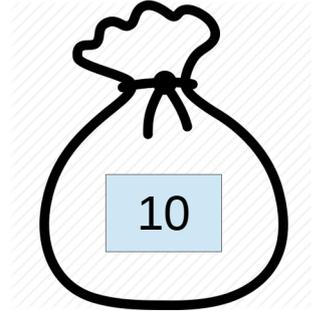- Let B ~ Unif(FIT). That is,
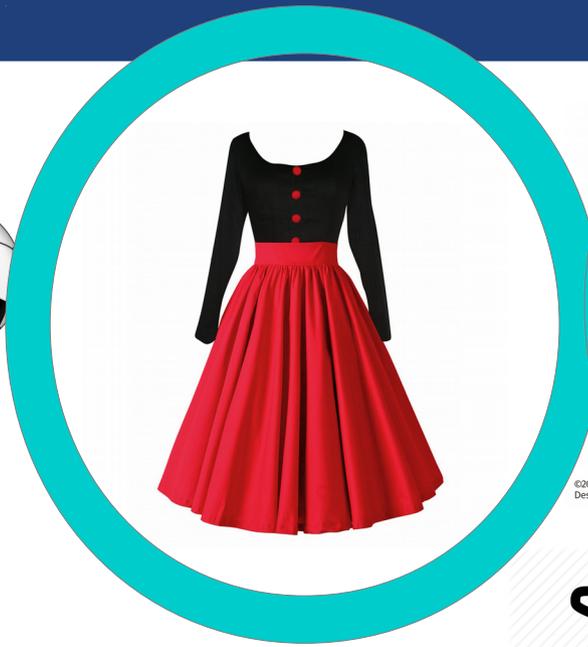
$$P(B = A) = 1/|\mathsf{FIT}|$$

Estimate

$$\mathbb{E}\left[\sum_{b \in B} value(b)\right].$$
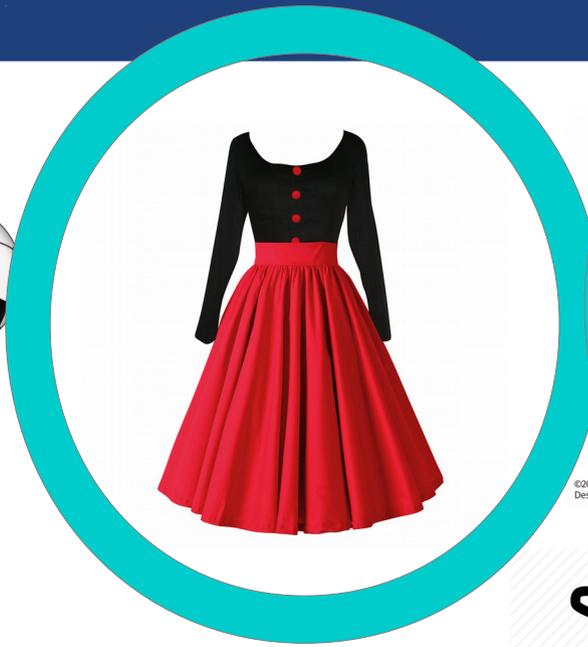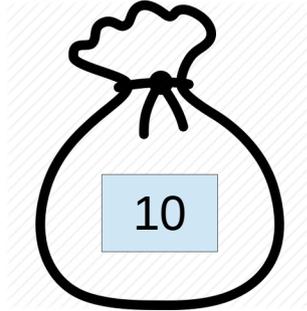
# How about rejection sampling?

10

10

10

10

10
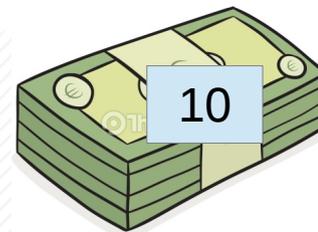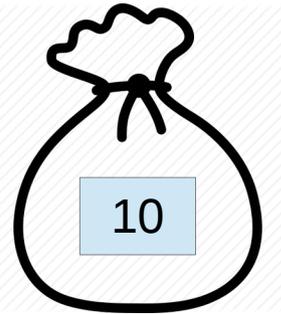
©2016 Sony Interactive Entertainment Inc. All rights reserved.
Design and specifications are subject to change without notice.

10

10

840

10

$$1/2\left(\frac{840}{\text{(dress + PS4)}} + \frac{10}{\text{(ball)}}\right)$$

# Let's implement it

# Generating samples from Q using MCMC

- Ingredients:
  - A prob. algo. T that transforms one sample into another.

  - A proof that T is "good" for Q.

- Recipe:
  - Take any sample x (not necessarily random).

  - For N sufficiently large, let x' = $T^N(x)$ = T(...(T(x))…).

  - Return x'

# What makes T "good" for Q?

- Let $\Omega = \{0,1,2,3\}$ and $Q = \text{Unif}(\Omega)$.



$(n - 1) \% |\Omega|$

0.25

0.5

n

0.25

$(n + 1) \% |\Omega|$

n → T

# What makes T "good" for Q?

- Let $\Omega = \{0,1,2,3\}$ and $Q = \text{Unif}(\Omega)$.

# What makes T "good" for Q?

- Let $\Omega = \{0,1,2,3\}$ and $Q = \text{Unif}(\Omega)$.



$n \rightarrow \boxed{T}$

- $0.25 \rightarrow (n - 1) \% |\Omega|$
- $0.5 \rightarrow n$
- $0.25 \rightarrow (n + 1) \% |\Omega|$

T induces an ergodic Markov chain!

# What makes T "good" for Q?

- Let $\Omega = \{0,1,2,3\}$ and $Q = \text{Unif}(\Omega)$.



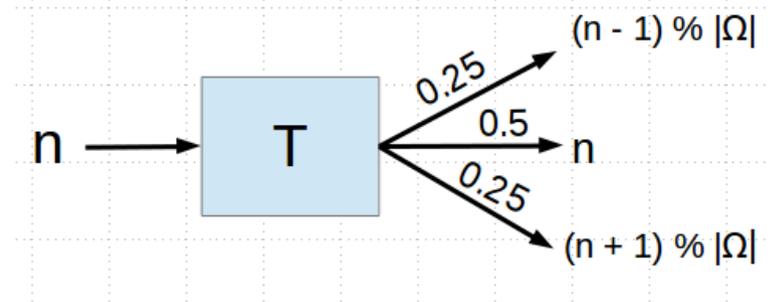T induces an ergodic Markov chain! Any state can reach any other state In 2,000,000 steps.

# What makes T "good" for Q?

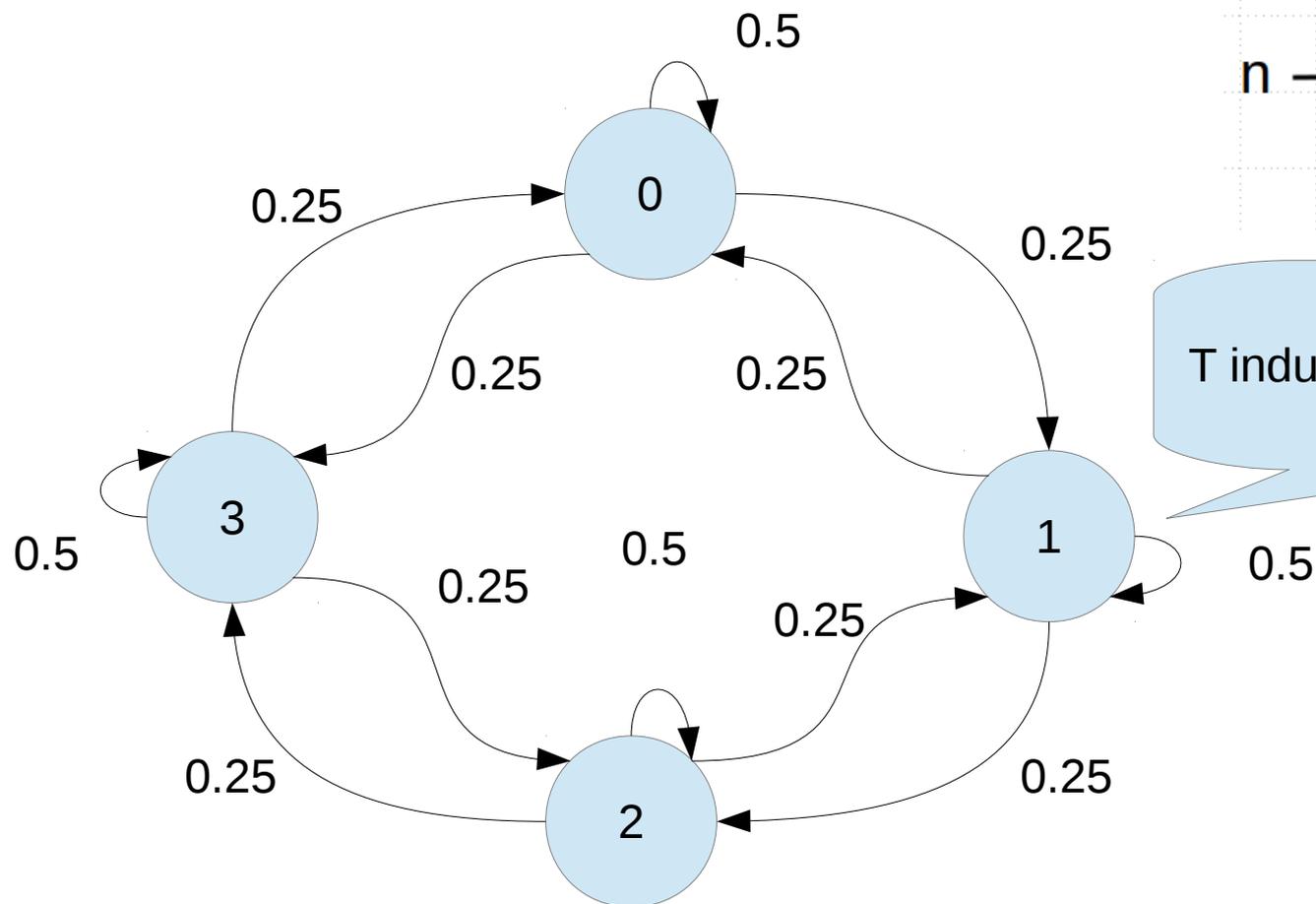- Let $\Omega = \{0,1,2,3\}$ and $Q = \text{Unif}(\Omega)$.



$n \longrightarrow \boxed{T} \begin{array}{l} \overset{0.25}{\longrightarrow} (n-1) \% |\Omega| \\ \overset{0.5}{\longrightarrow} n \\ \overset{0.25}{\longrightarrow} (n+1) \% |\Omega| \end{array}$

T induces an ergodic Markov chain! Any state can reach any other state In 2,000,000 steps.

$Q(x)\, R(x \mid x') = Q(x')\, R(x' \mid x)$, for any two states x and x'.

Markov graph of T

# What makes T "good" for Q?

- Let **Ω = {0,…, 9999}** and Q = Unif(Ω).

Markov graph of T

# What makes T "good" for Q?

- Let M be the Markov chain induced by T and let R be M's transition probability.
- T is "good" for Q if
  - M is ergodic.
  - $Q(x)R(x'|x) = Q(x')R(x'|x)$, for all $x$, $x'$.

# What makes T "good" for Q?

- Let M be the Markov chain induced by T and let R be M's transition probability.
- T is "good" for Q if
  - M is ergodic.
  - $Q(x)R(x'|x) = Q(x')R(x'|x)$, for all $x$, $x'$.

Warning, these are sufficiency conditions!
There may be other algorithms that are
"good" for Q, but do not satisfy these
conditions.

# Computing expected loss from a burglary

# Generating samples from a complex distribution Q using MCMC

- Ingredients:
  - A prob. algo. T that transforms one sample into another.

  - A proof that T is "good" for Q.

- Recipe:
  - Take any sample x (not necessarily random).

  - For N sufficiently large, let x' = $T^N(x)$ = T(...(T(x))…).

  - Return x'

# A "good" T for
# the uniform distr. on $2^{Col}$

- Let B in FIT.
  - Flip a coin. If heads, then return B.
  - Pick an object b in Col uniformly at random.
  - If b in B:
    - return B \ {b}
  - If b not in B:
    - If the total weight of B U {b} <= Cap:
      - return B U {b}
    - Else:
      - return B

# Proving T is "good" for Q

- Part 1 of 2. Show T induces an ergodic Markov chain.

# Proving T is "good" for Q

- Part 1 of 2. Show T induces an ergodic Markov chain.
  - Insight 1: The Markov graph of T is connected.
    - If you are lucky enough, T transforms any B into the empty set after some steps. If you are even luckier, then T transforms the empty set into B' after some steps.

# Proving T is "good" for Q

- Part 1 of 2. Show T induces an ergodic Markov chain.
    - Insight 1: The Markov graph of T is connected.
        - If you are lucky enough, T transforms any B into the empty set after some steps. If you are even luckier, then T transforms the empty set into B' after some steps.
    - Insight 2: If L >> $|2^{Col}|$, then T can reach any B' from any B in at most L steps.

# Proving T is "good" for Q

- Part 1 of 2. Show T induces an ergodic Markov chain.
  - Insight 1: The Markov graph of T is connected.
    - If you are lucky enough, T transforms any B into the empty set after some steps. If you are even luckier, then T transforms the empty set into B' after some steps.
  - Insight 2: If $L >> |2^{Col}|$, then T can reach any B' from any B in at most L steps.
  - Insight 3: There is a self-loop for every B' in the Markov graph.
    - If you happen to arrive to B' before t steps, just use the extra steps on the self-loop to arrive in exactly t steps.

# Proving T is "good" for Q

- Part 2 of 2. $Q(x)\, R(x \mid x') = Q(x')\, R(x' \mid x)$, for any $x$, $x'$.

# Proving T is "good" for Q

- Part 2 of 2. $Q(x) R(x \mid x') = Q(x') R(x' \mid x)$, for any $x, x'$.
  - Hint 1: $Q(x) = Q(x')$.

# Proving T is "good" for Q

- Part 2 of 2. Q(x) R(x | x') = Q(x') R(x' | x), for any x, x'.
    - Hint 1: Q(x) = Q(x').
    - Hint 2: R(x' | x) = R(x' | x).