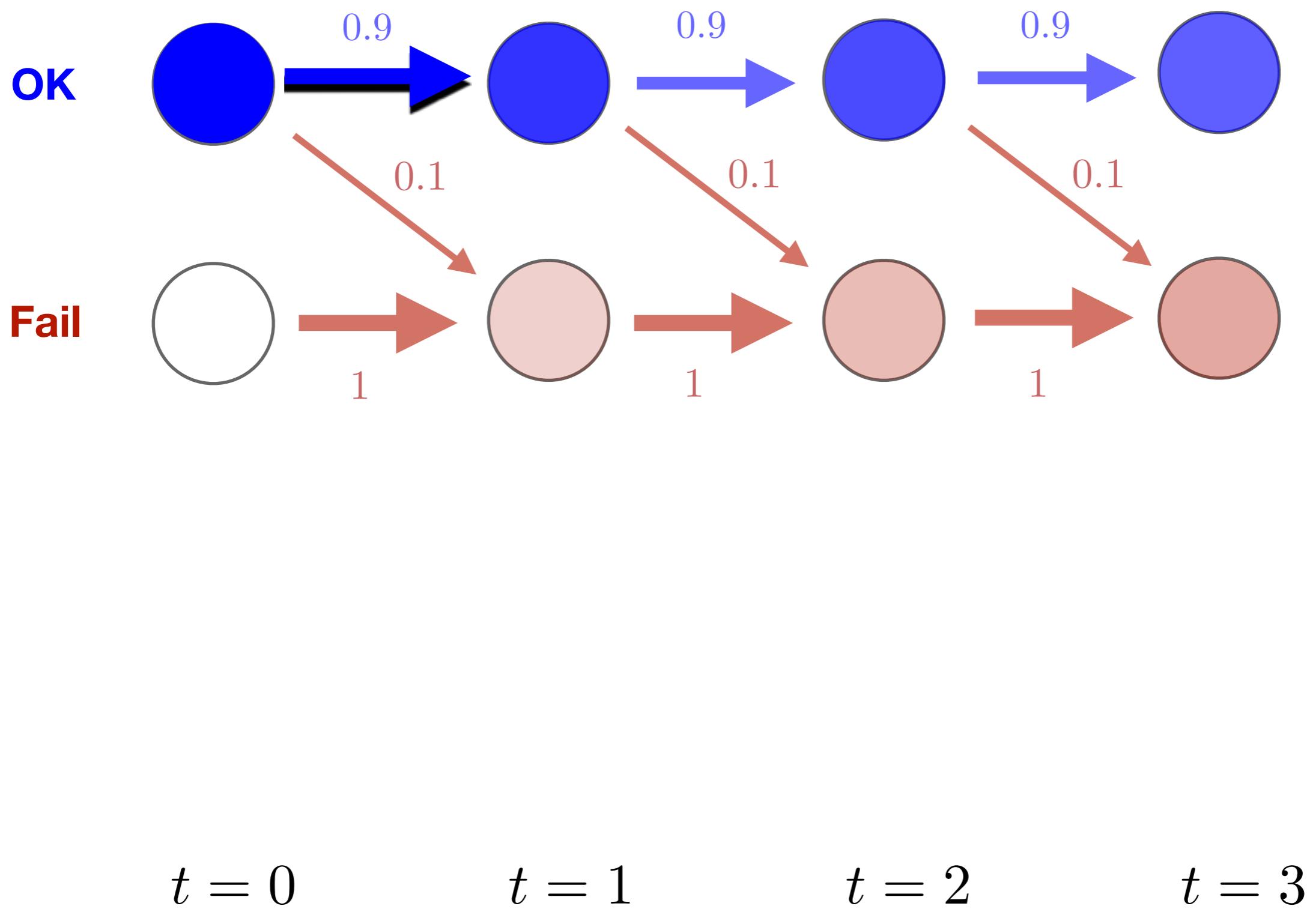
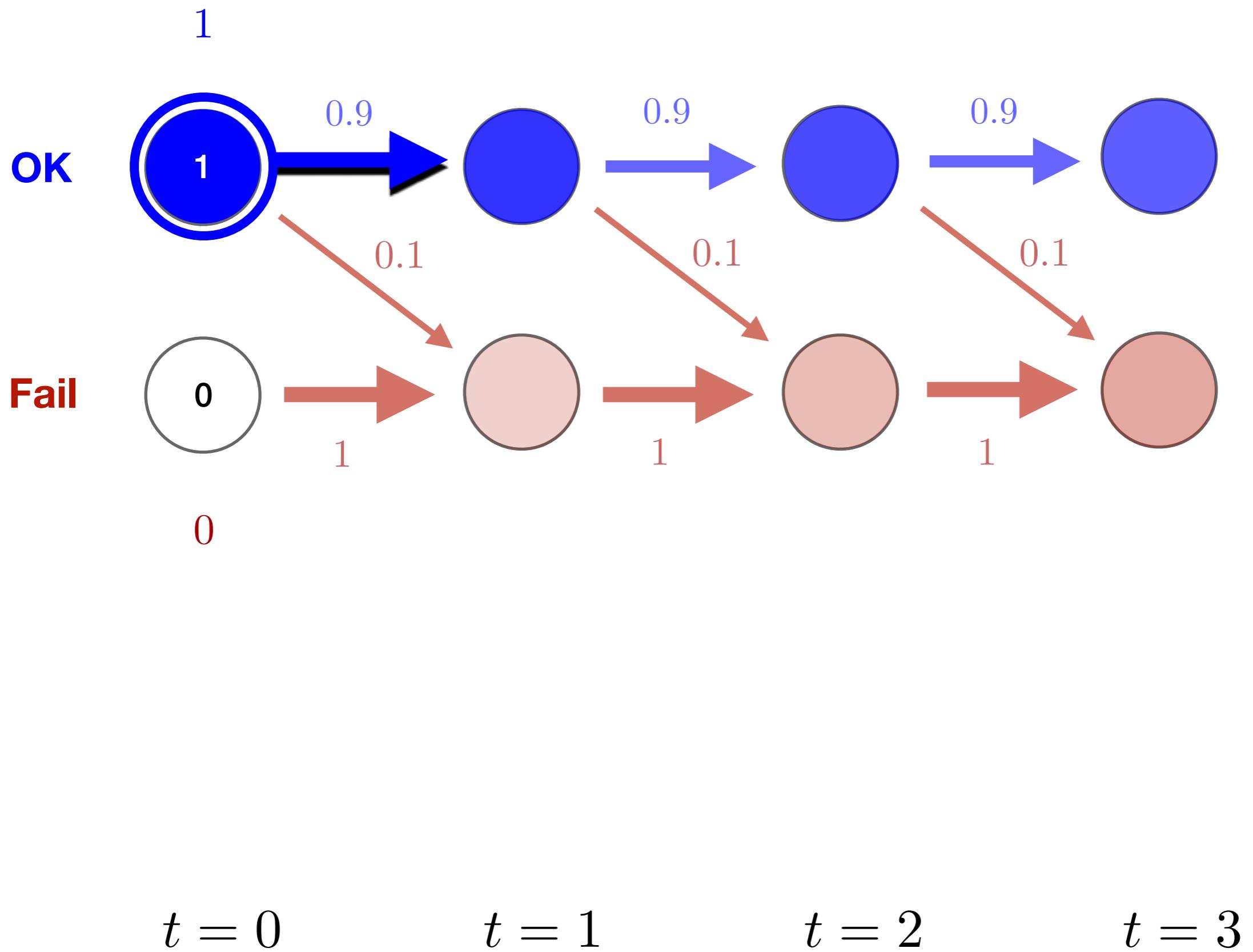


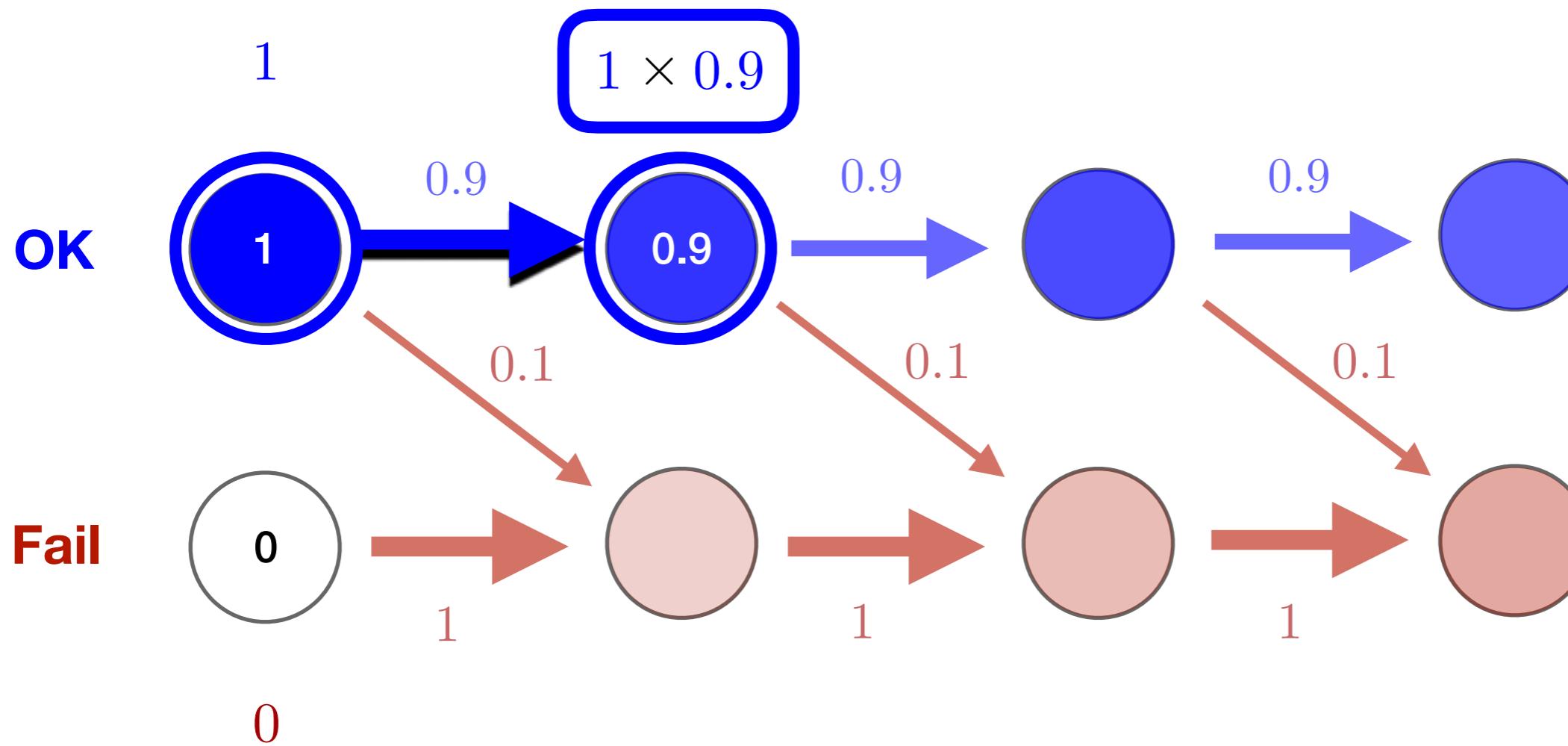
Markov Chain: Examples

Stochastic Processes

Major sources:





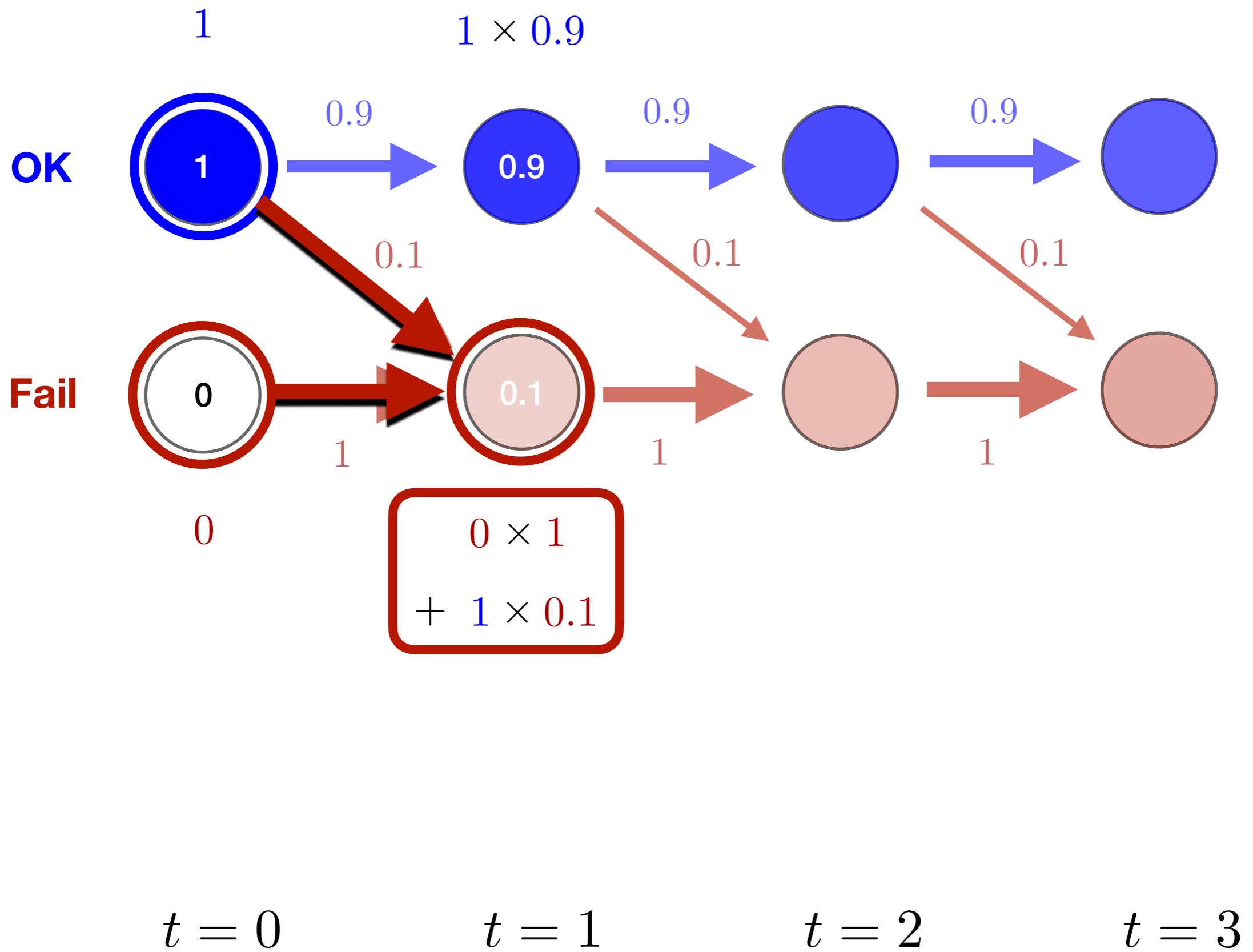


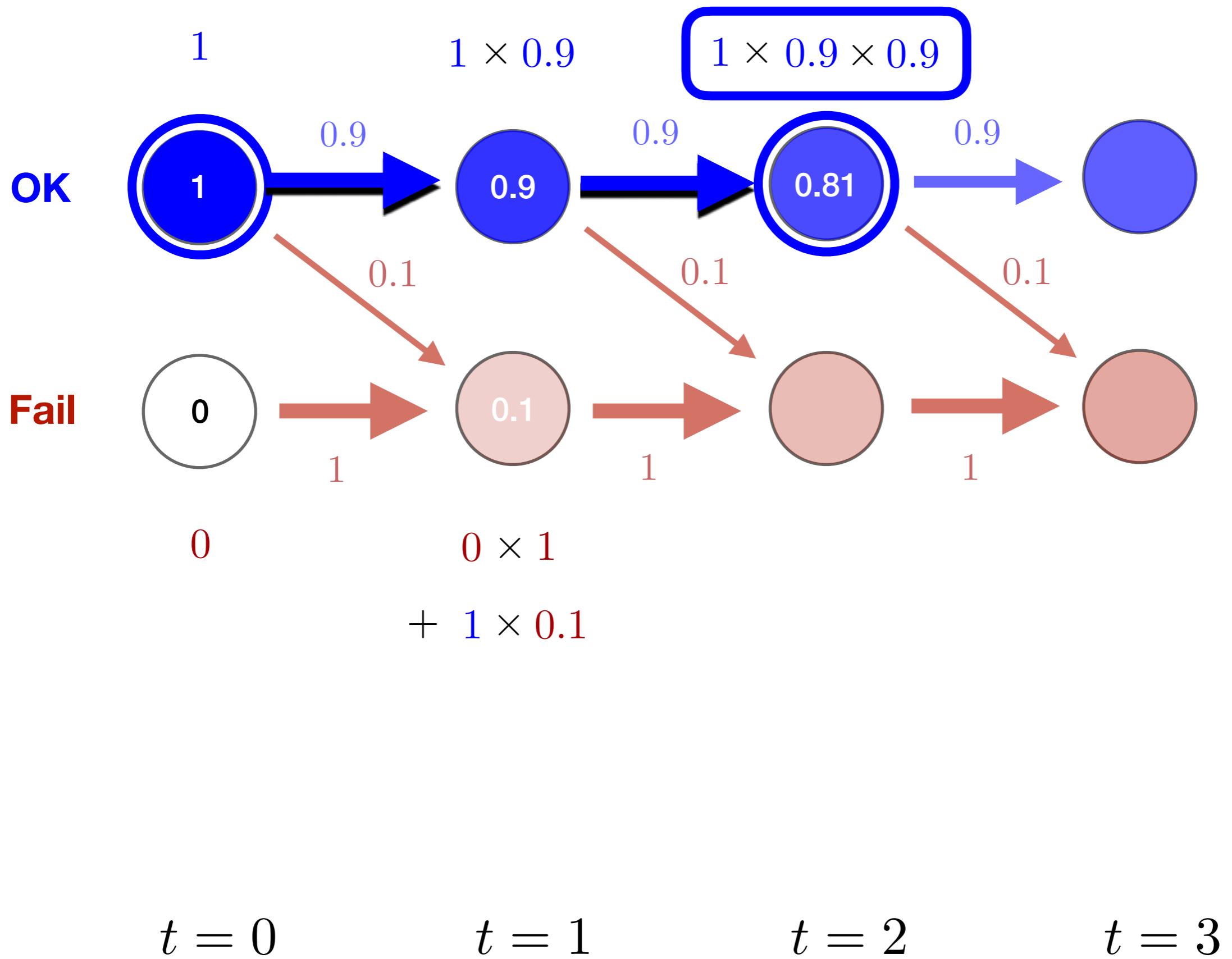
$$t = 0$$

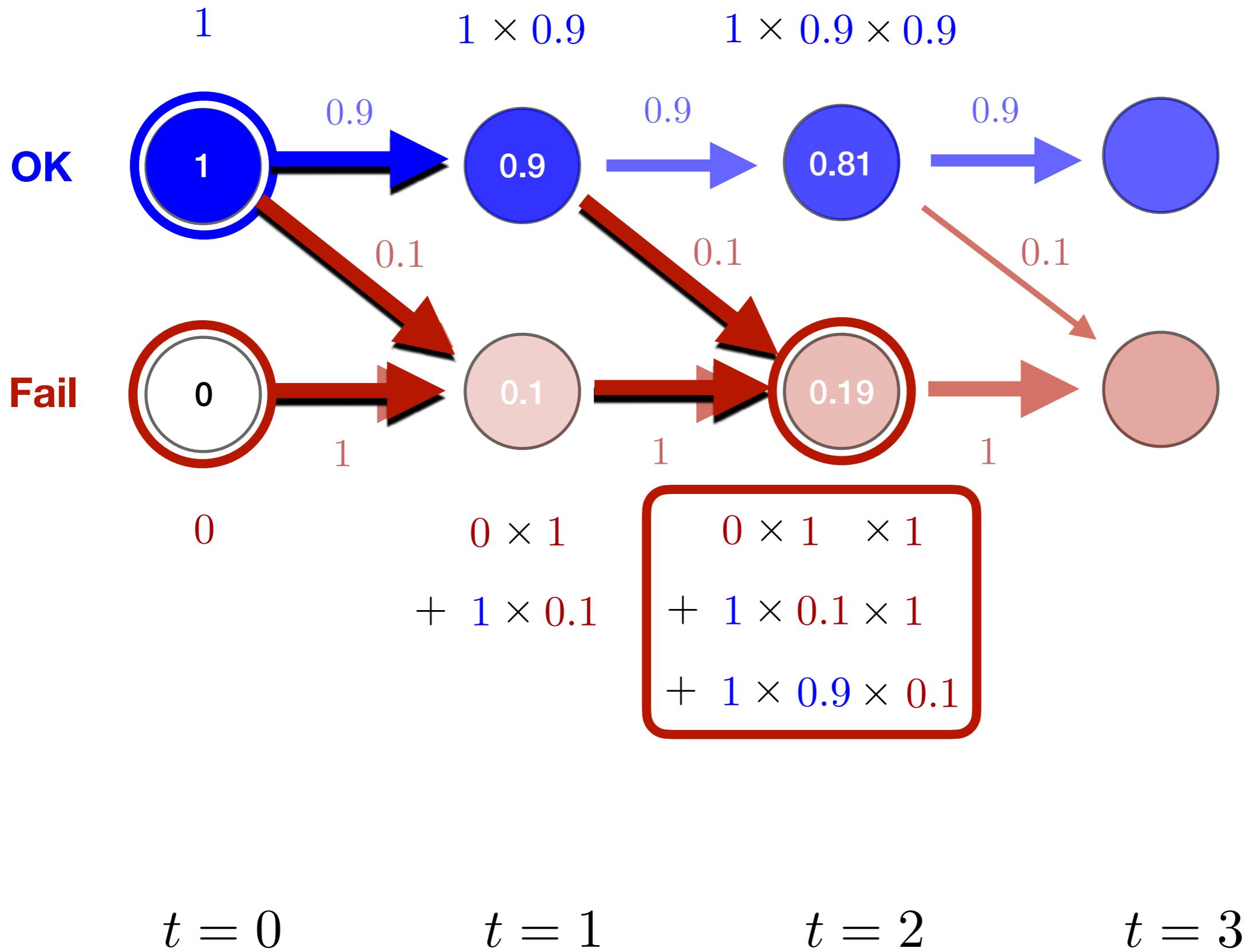
$$t = 1$$

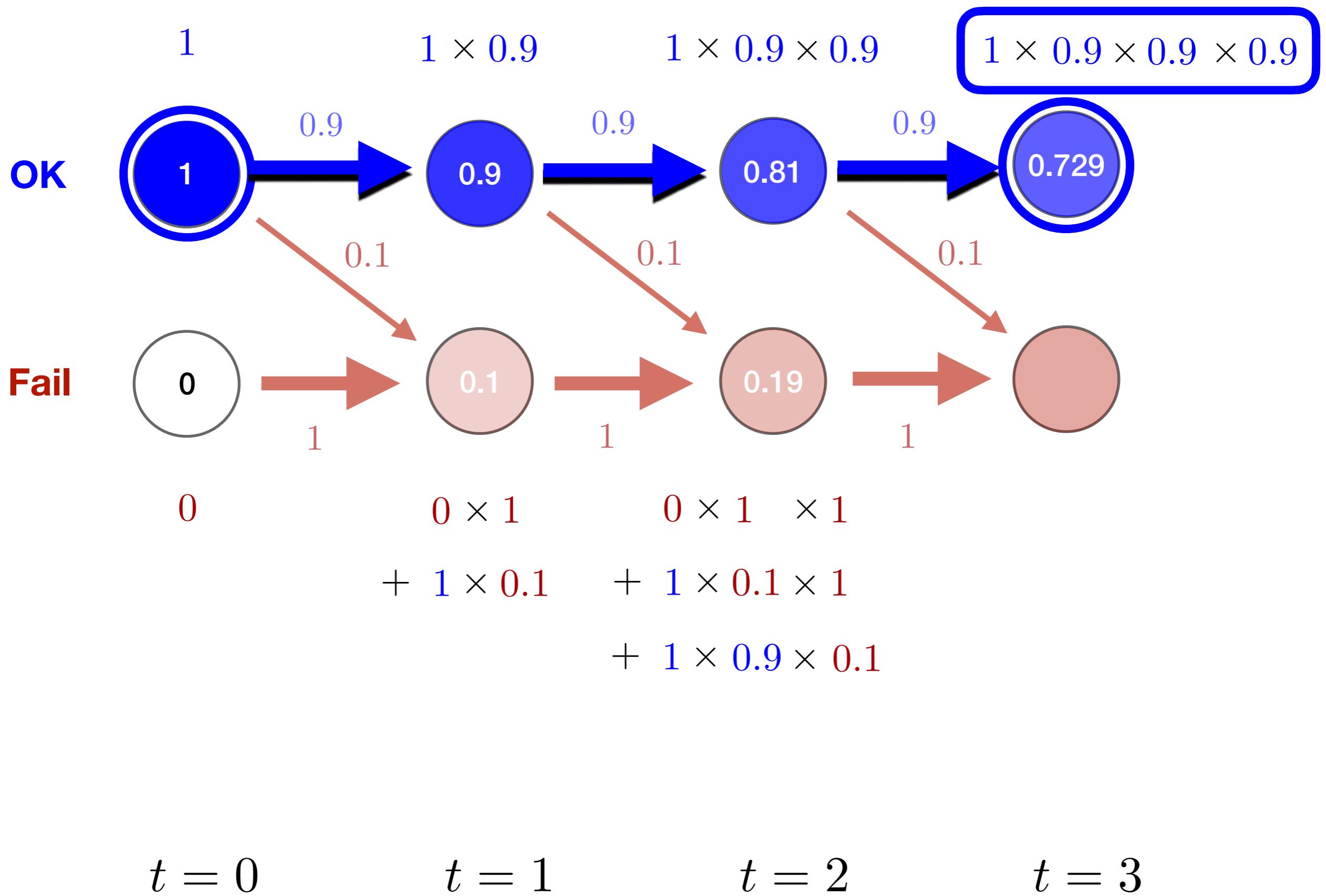
$$t = 2$$

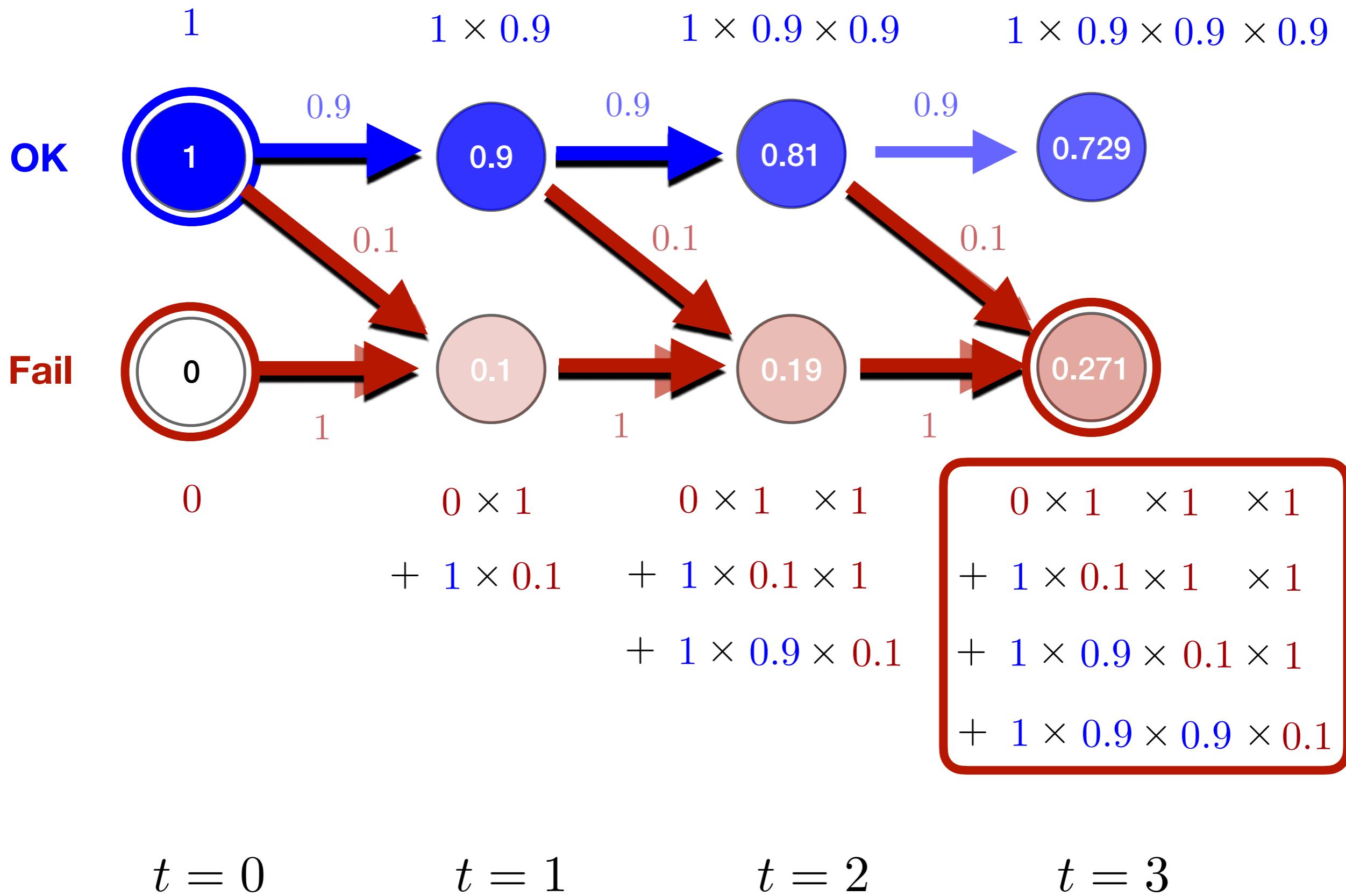
$$t = 3$$

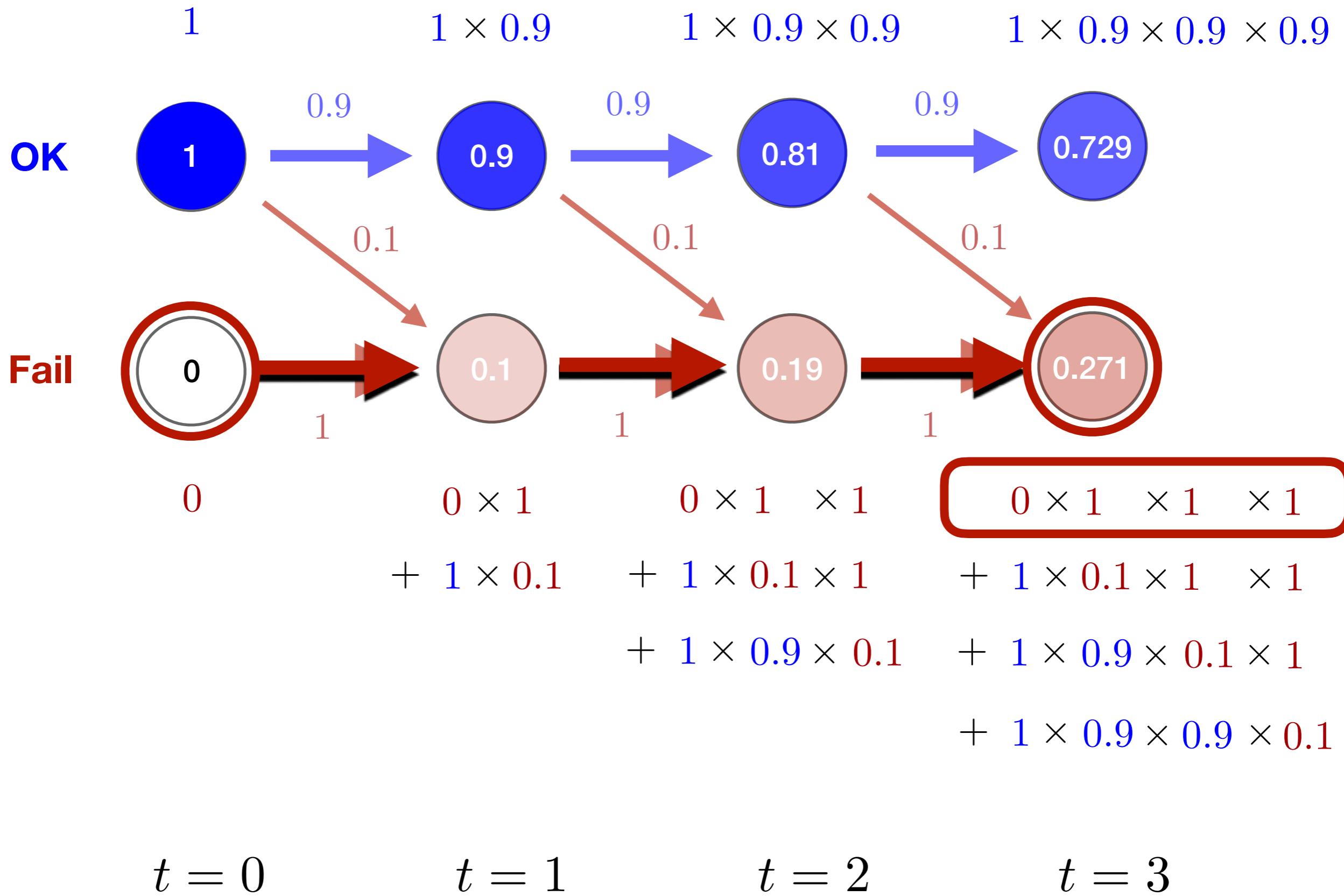


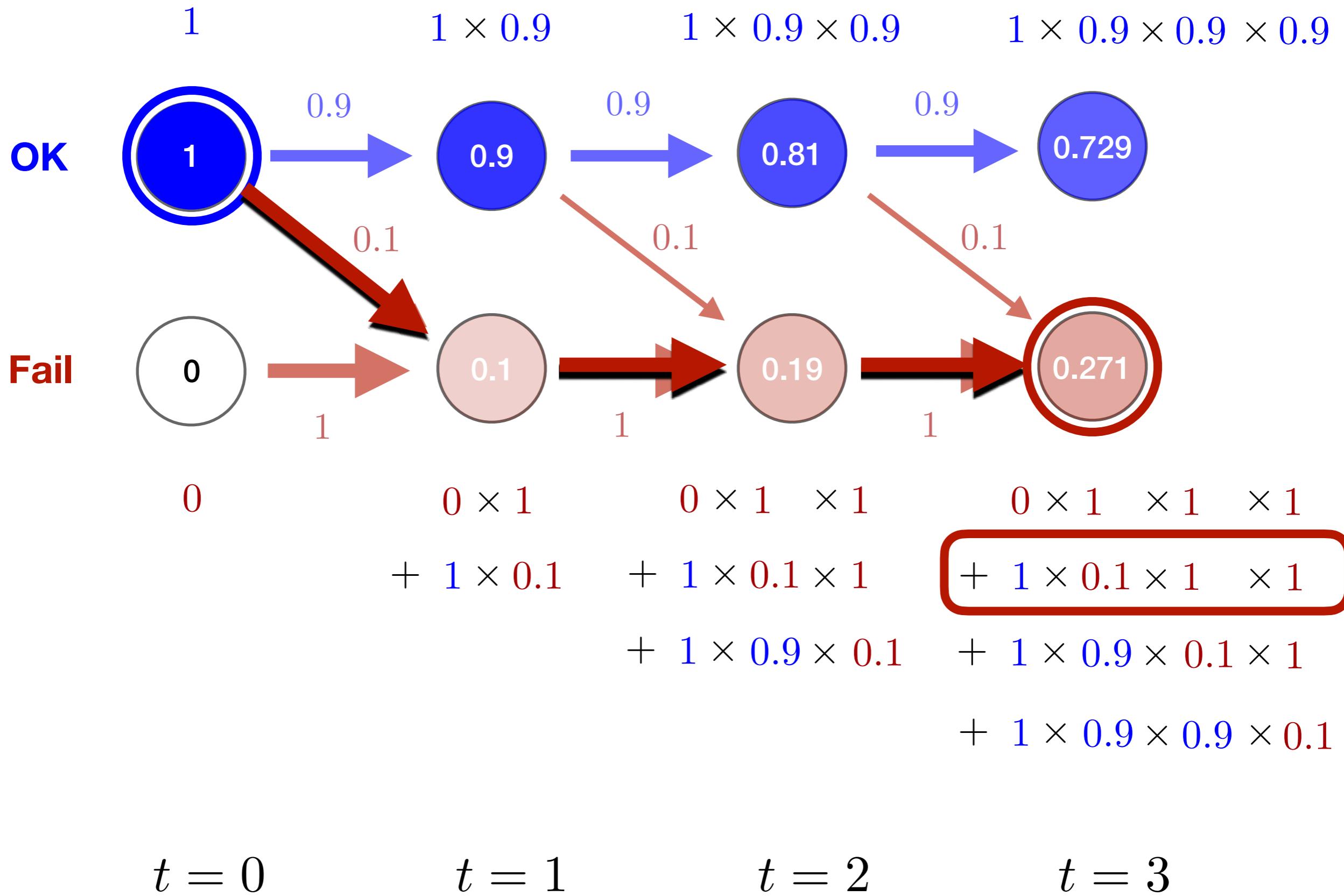


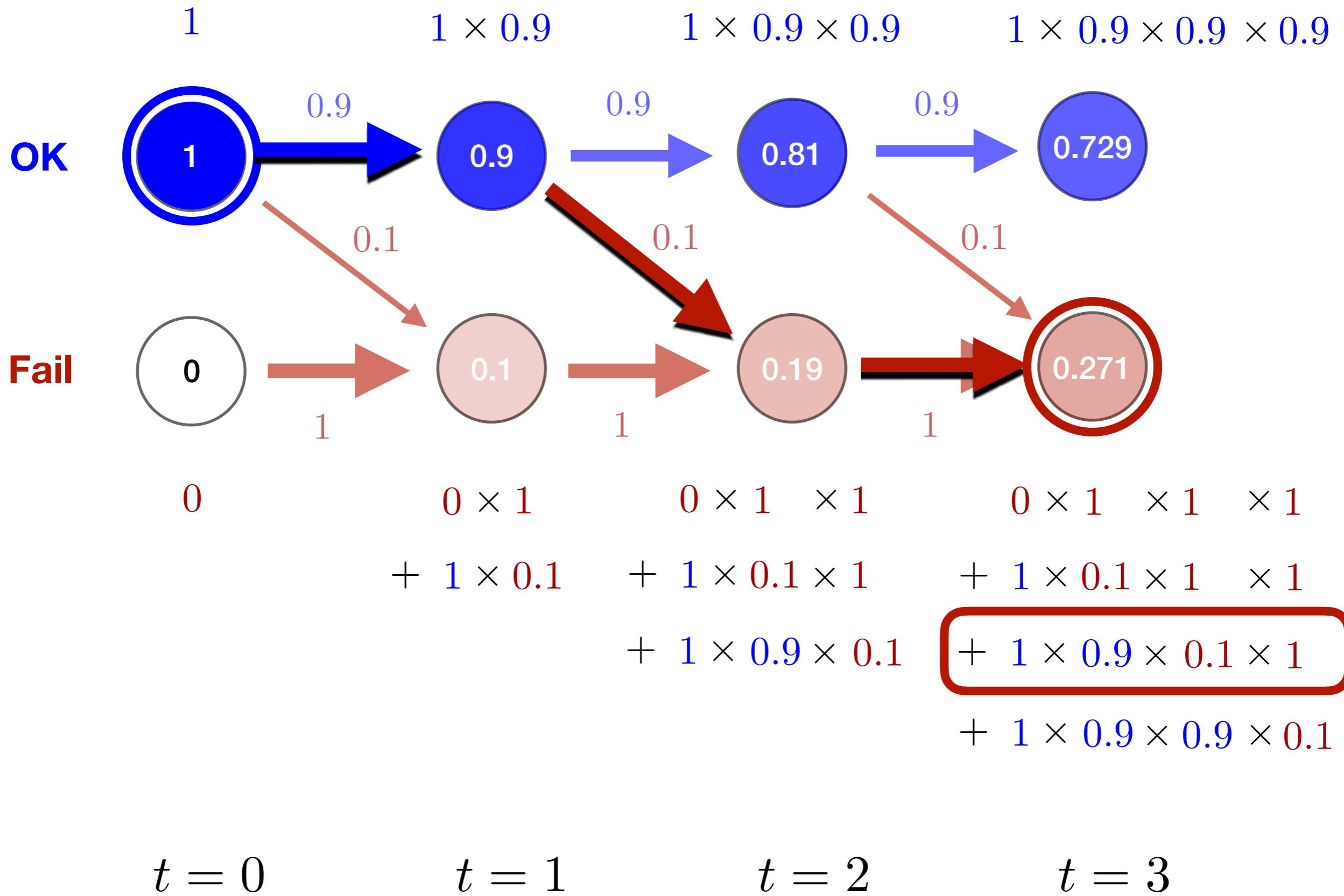


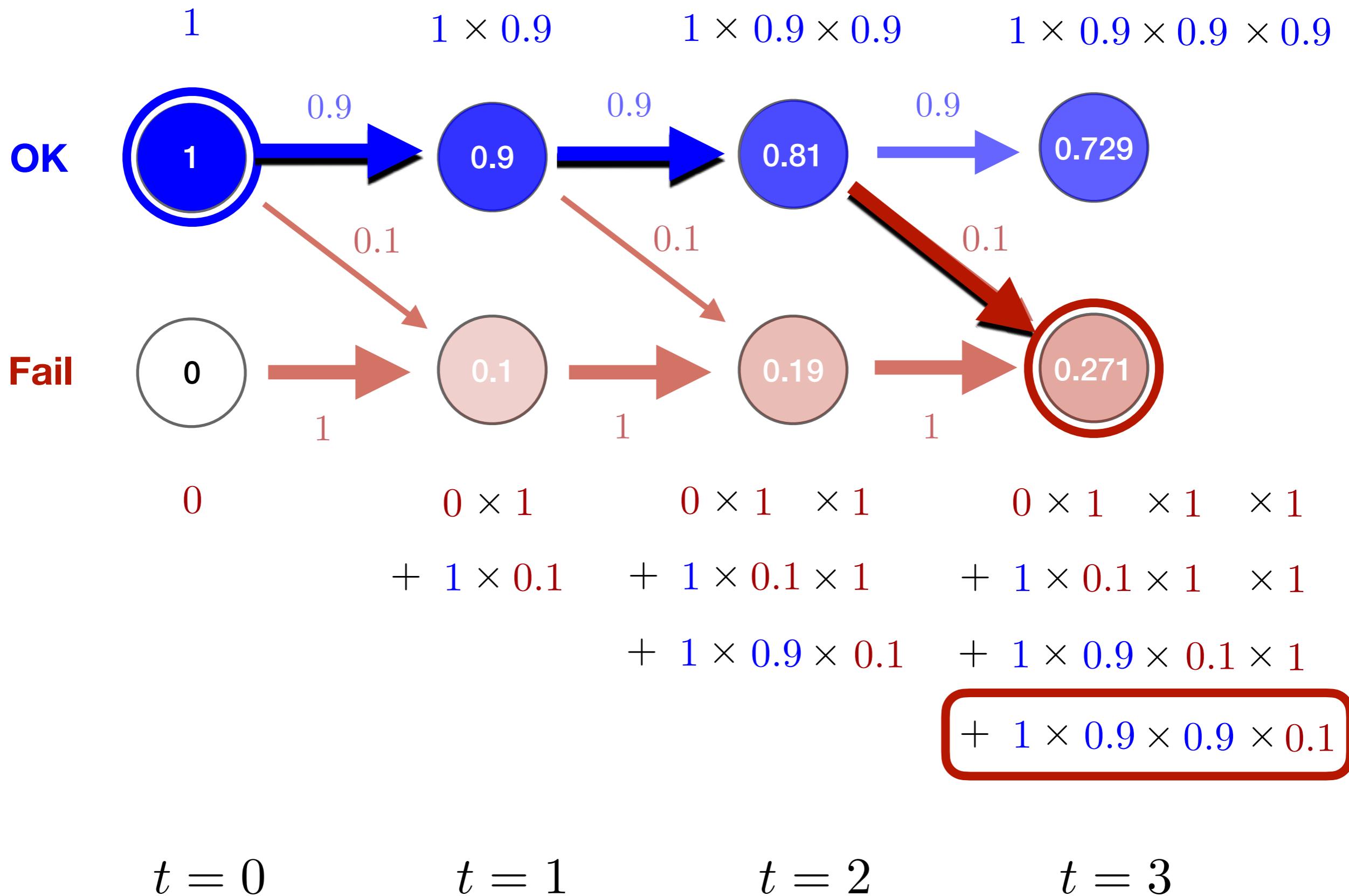


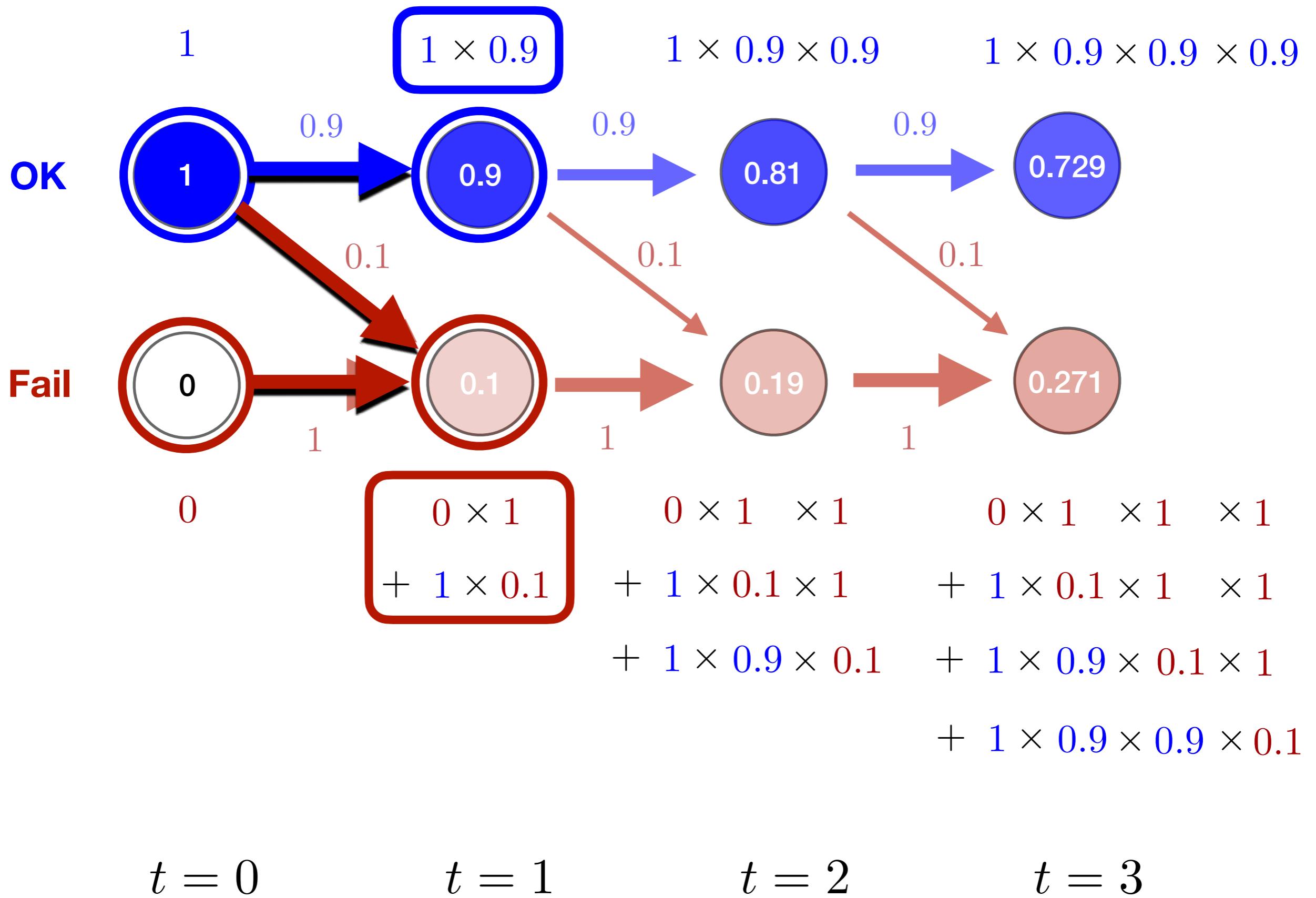


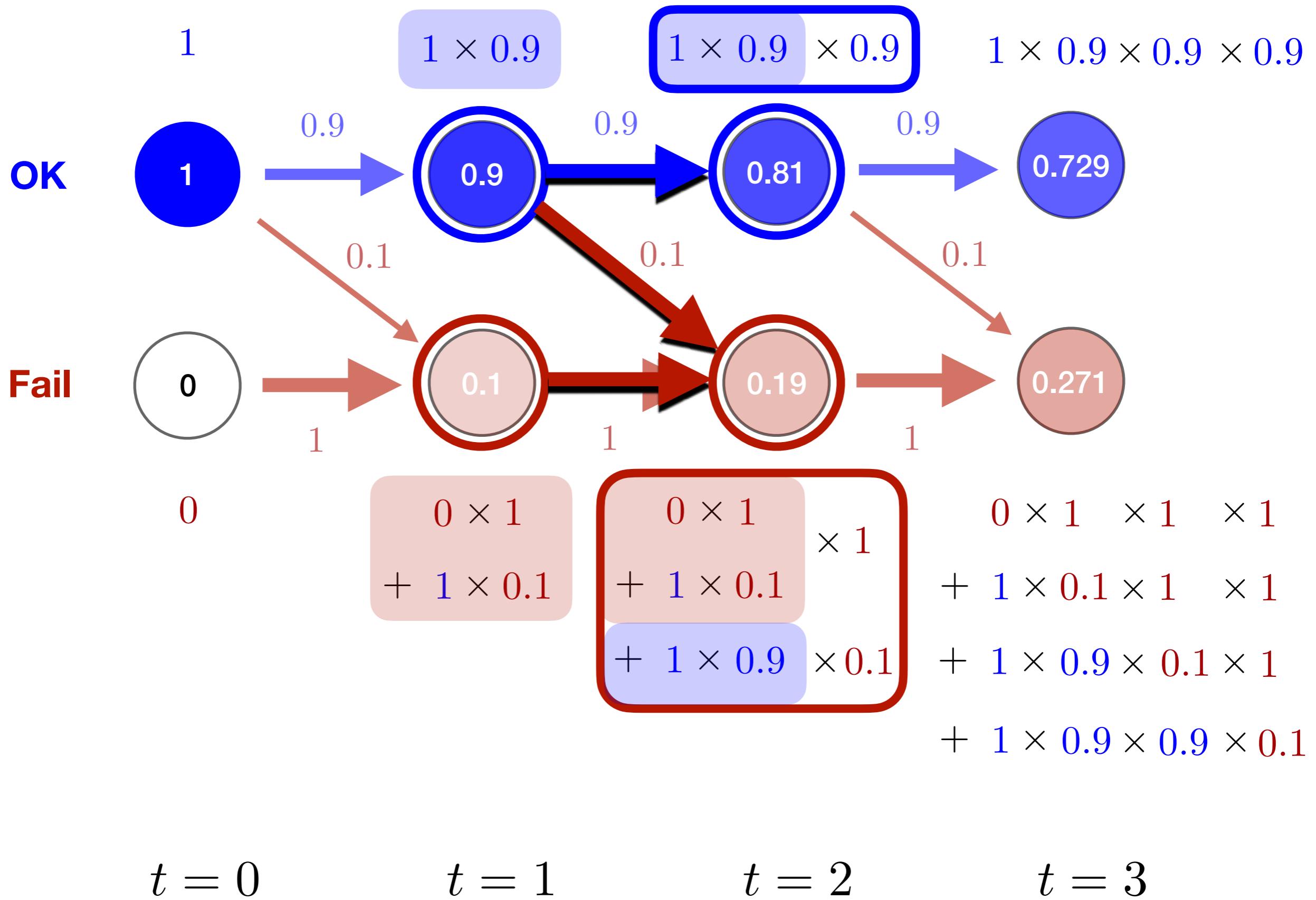


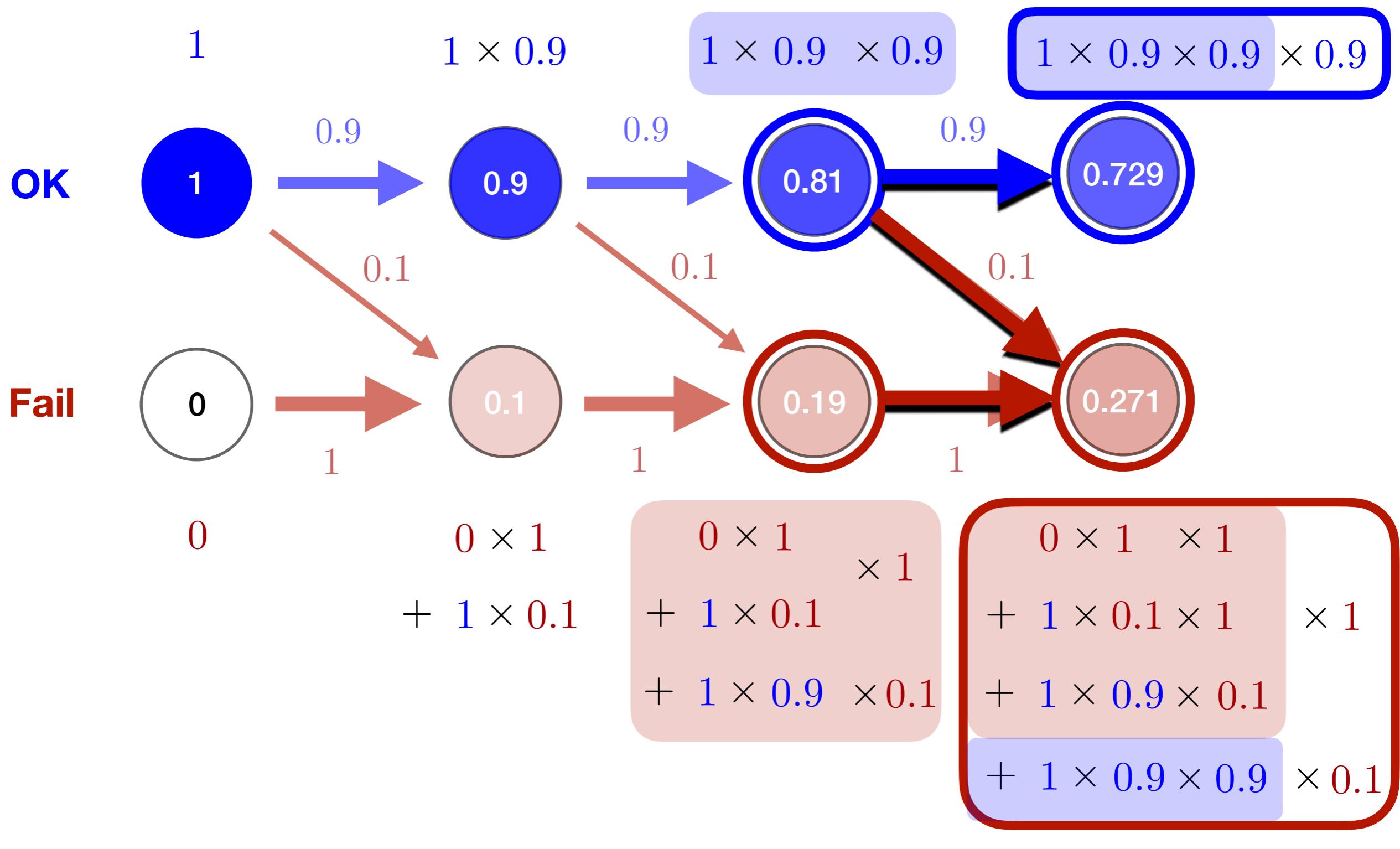


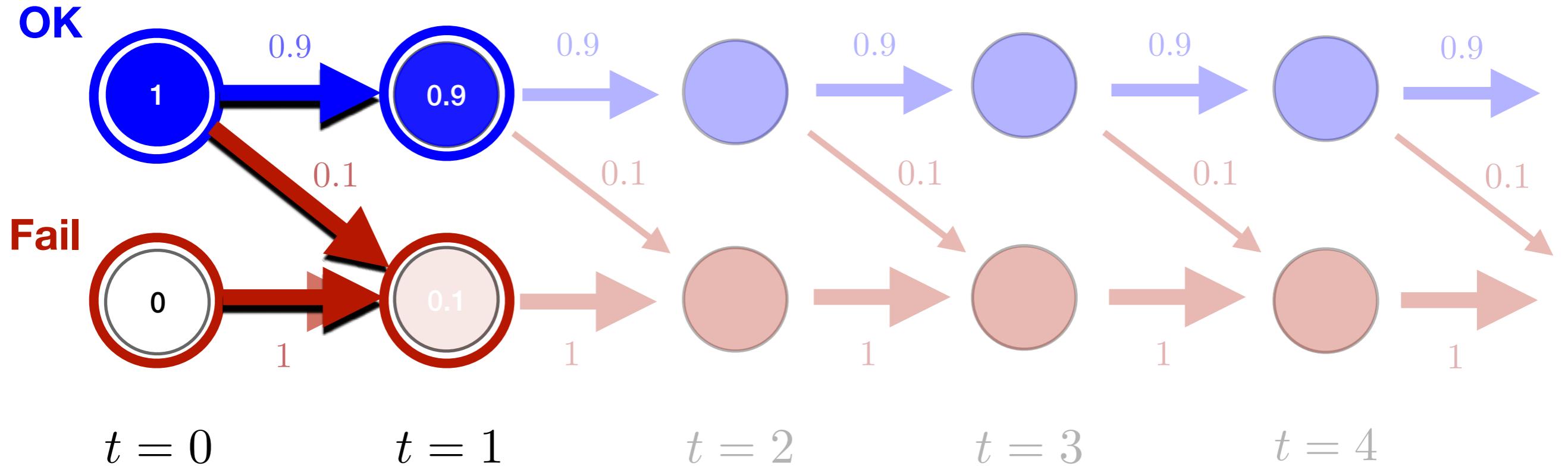




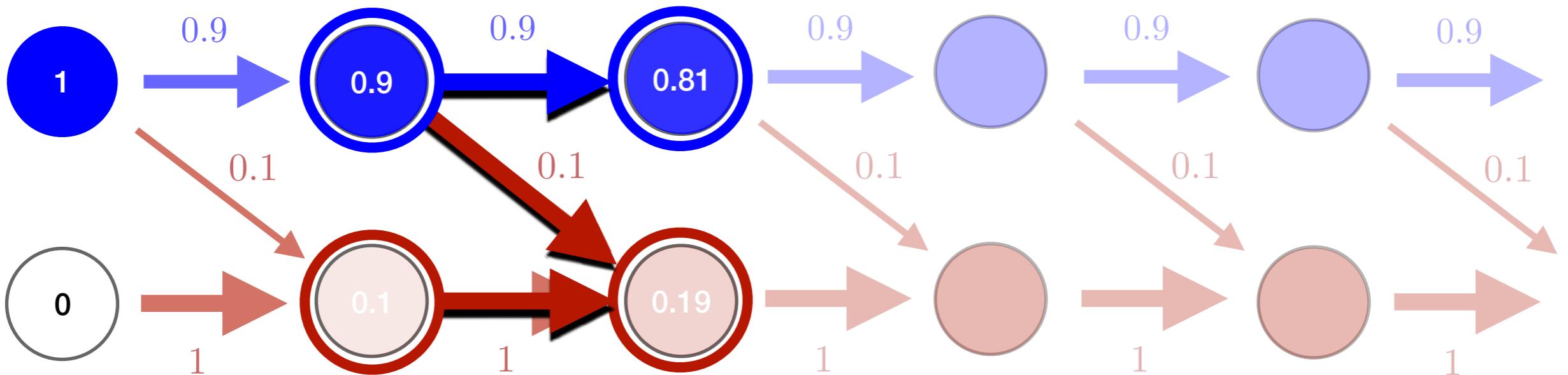








OK



$t = 0$

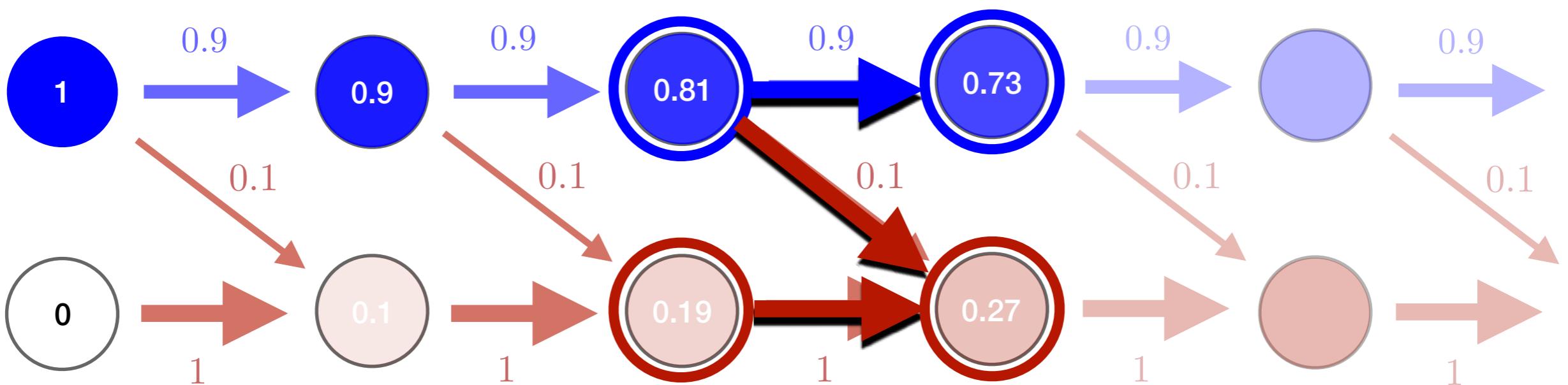
$t = 1$

$t = 2$

$t = 3$

$t = 4$

OK



$t = 0$

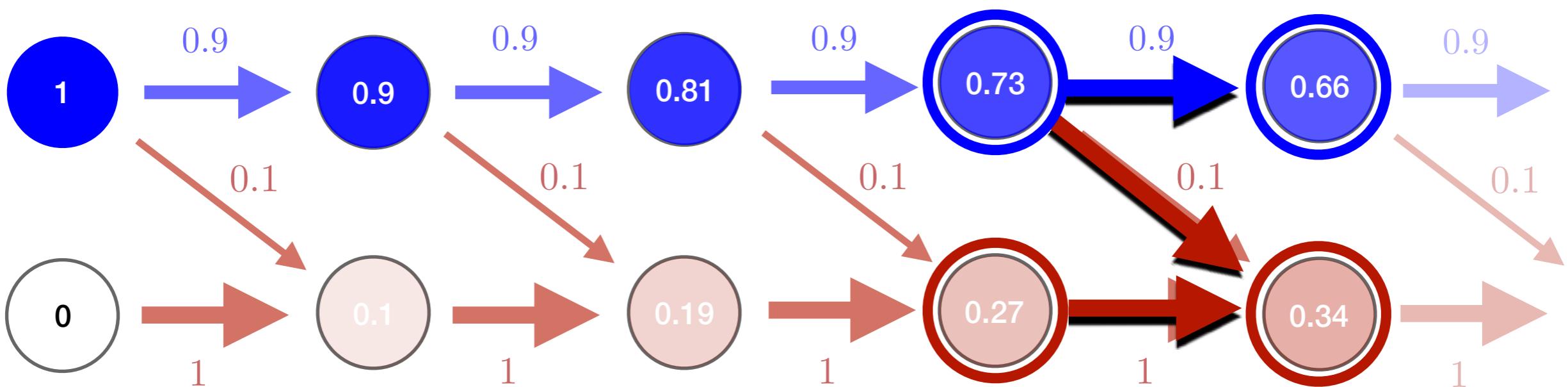
$t = 1$

$t = 2$

$t = 3$

$t = 4$

OK



$t = 0$

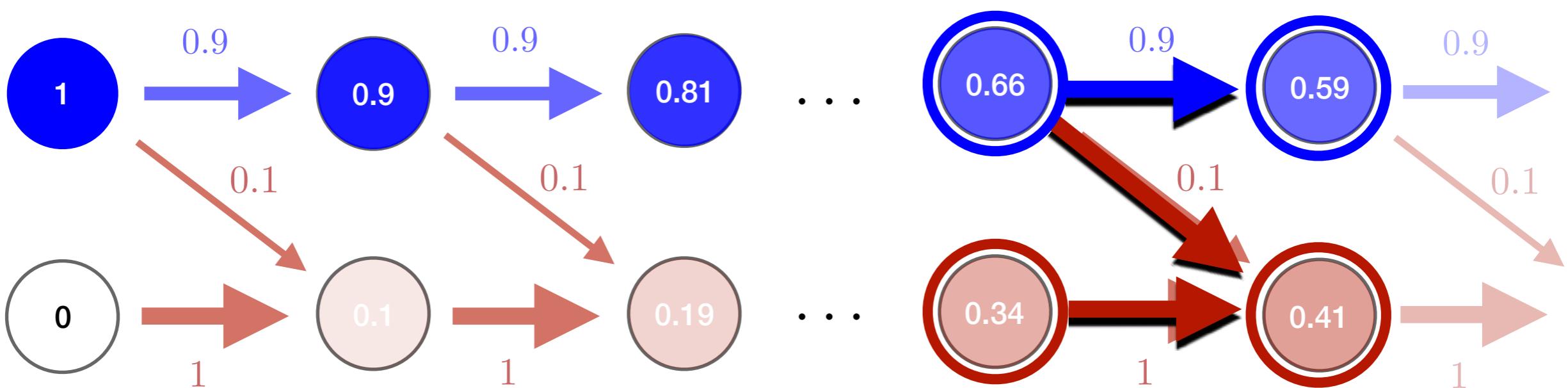
$t = 1$

$t = 2$

$t = 3$

$t = 4$

OK



$t = 0$

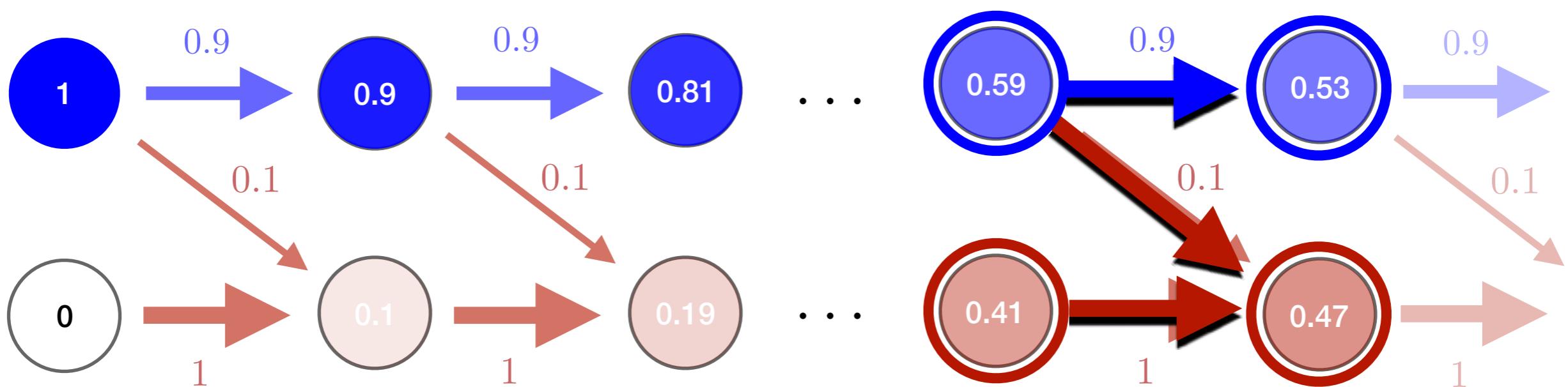
$t = 1$

$t = 2$

$t = 4$

$t = 5$

OK



$t = 0$

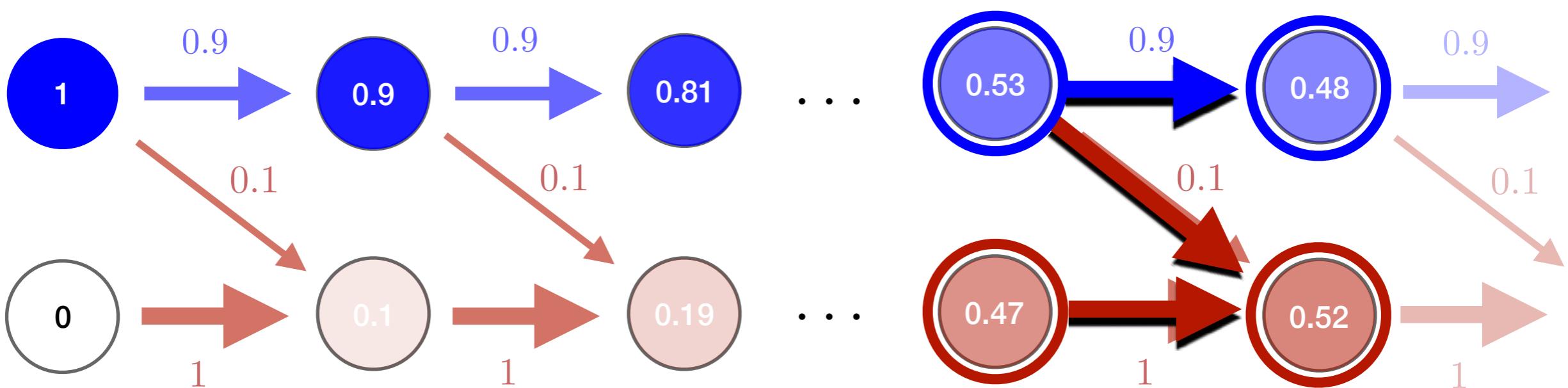
$t = 1$

$t = 2$

$t = 5$

$t = 6$

OK



$t = 0$

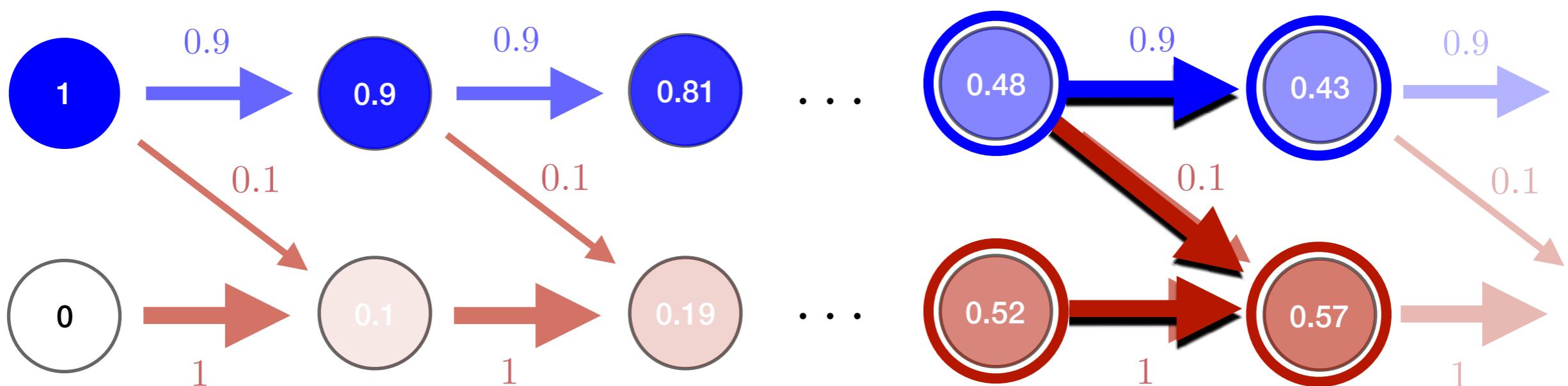
$t = 1$

$t = 2$

$t = 6$

$t = 7$

OK



$t = 0$

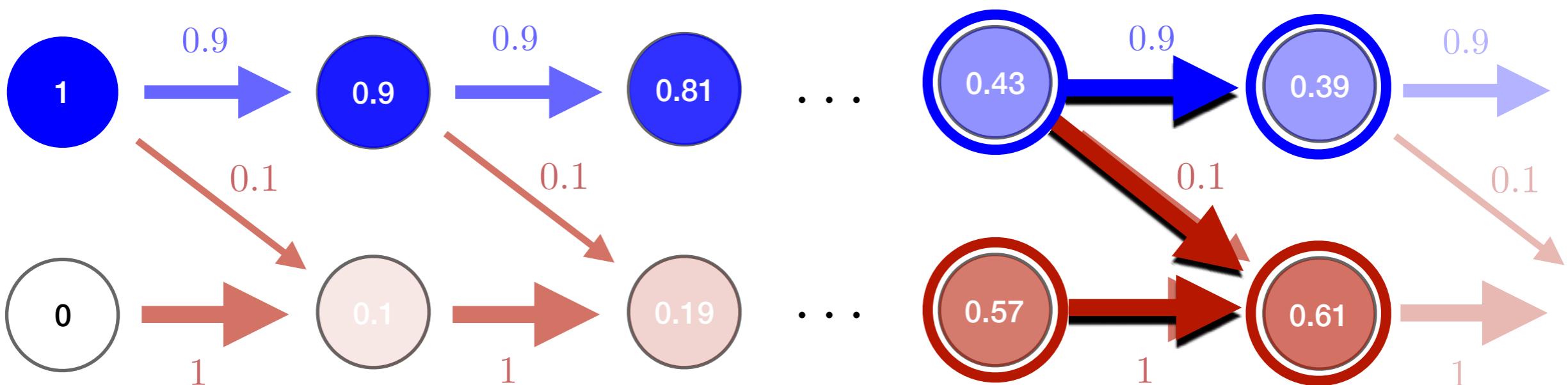
$t = 1$

$t = 2$

$t = 7$

$t = 8$

OK



$t = 0$

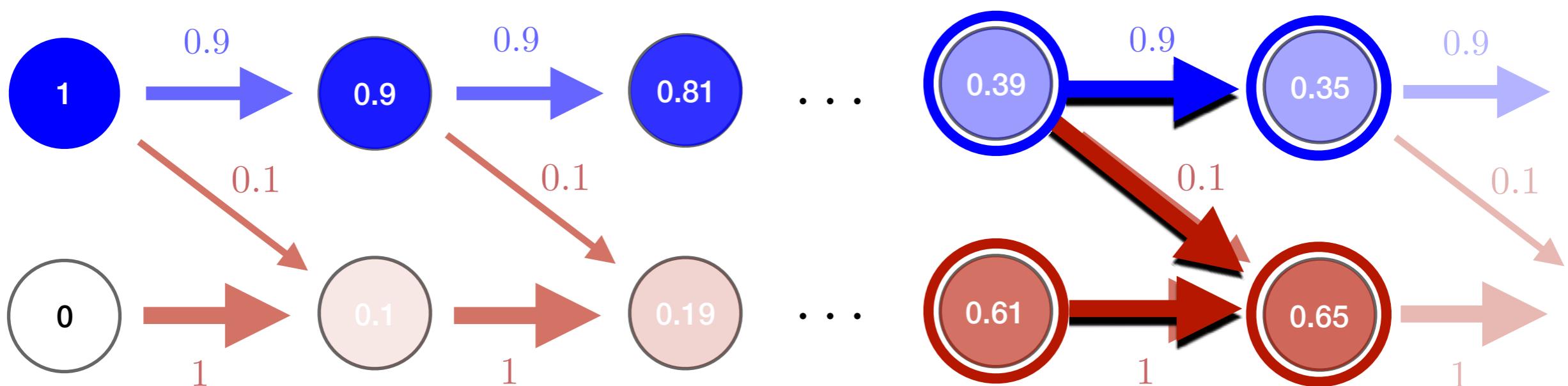
$t = 1$

$t = 2$

$t = 8$

$t = 9$

OK



$t = 0$

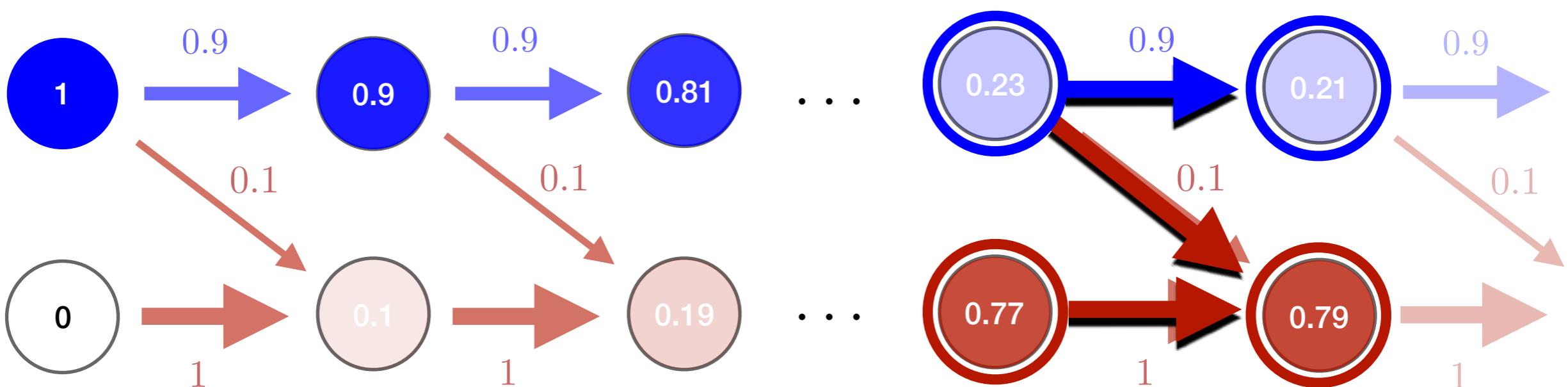
$t = 1$

$t = 2$

$t = 9$

$t = 10$

OK



$t = 0$

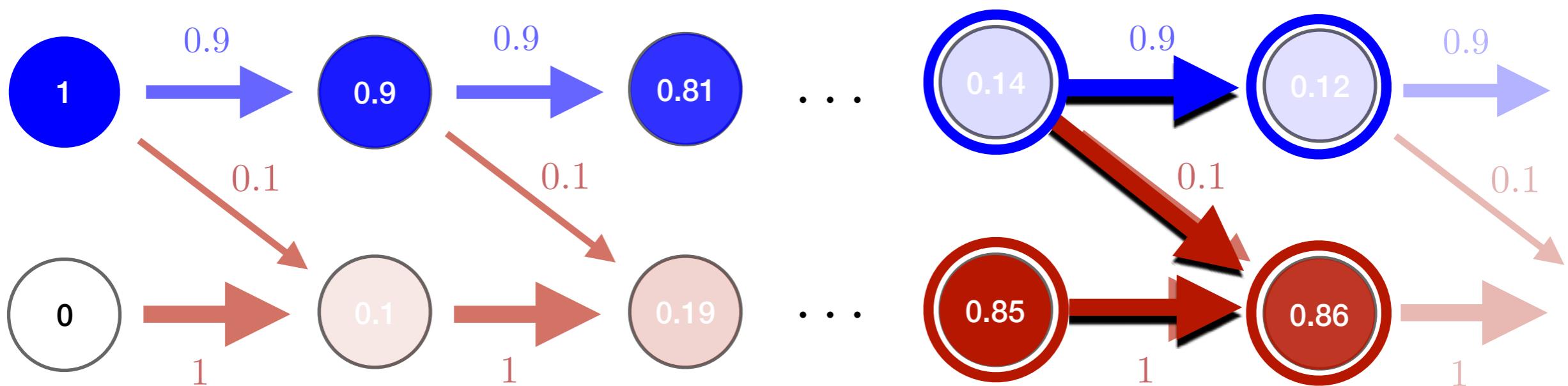
$t = 1$

$t = 2$

$t = 14$

$t = 15$

OK



$t = 0$

$t = 1$

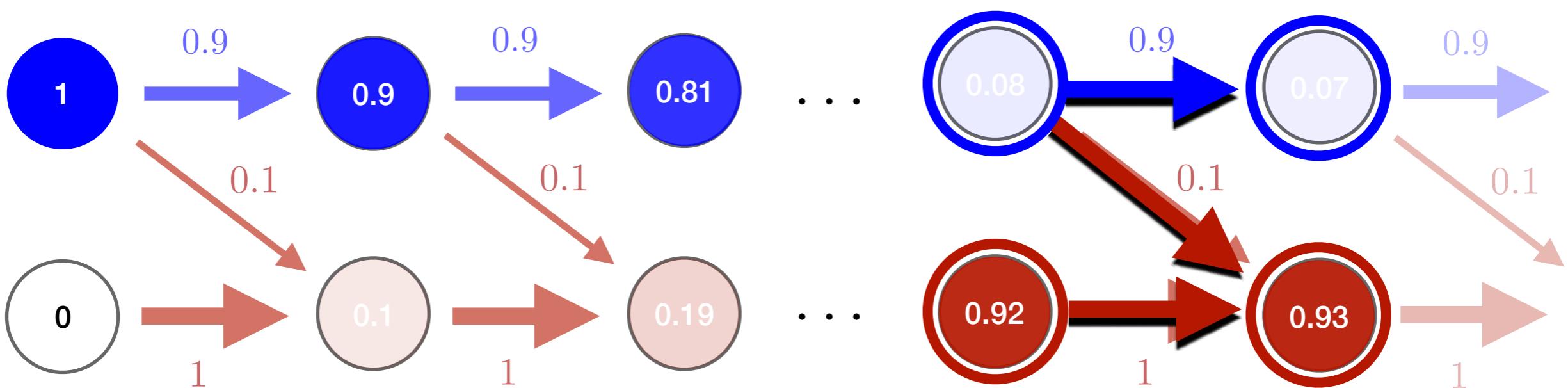
$t = 2$

$t = 19$

$t = 20$

Fail

OK



$t = 0$

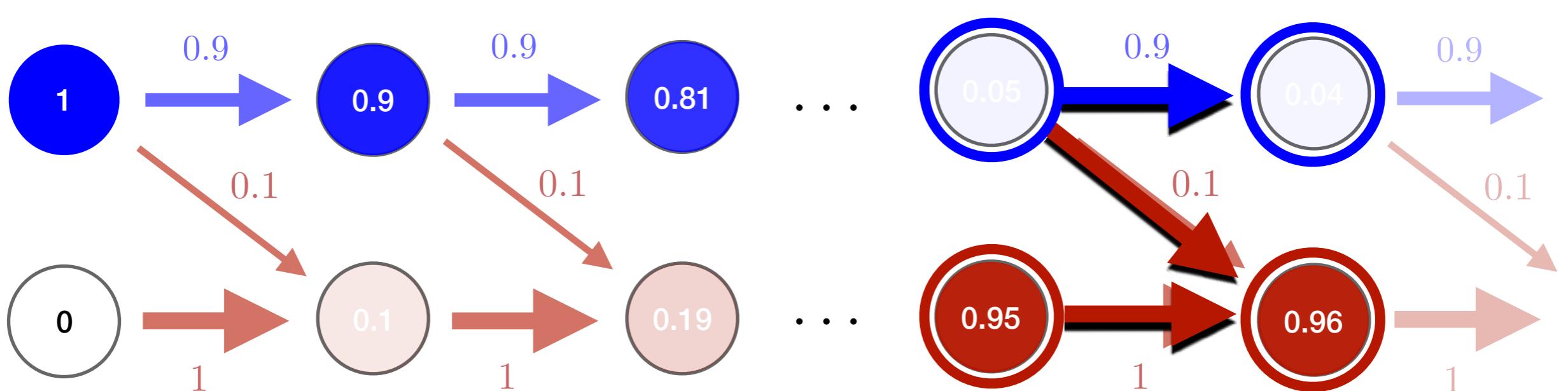
$t = 1$

$t = 2$

$t = 24$

$t = 25$

OK



$t = 0$

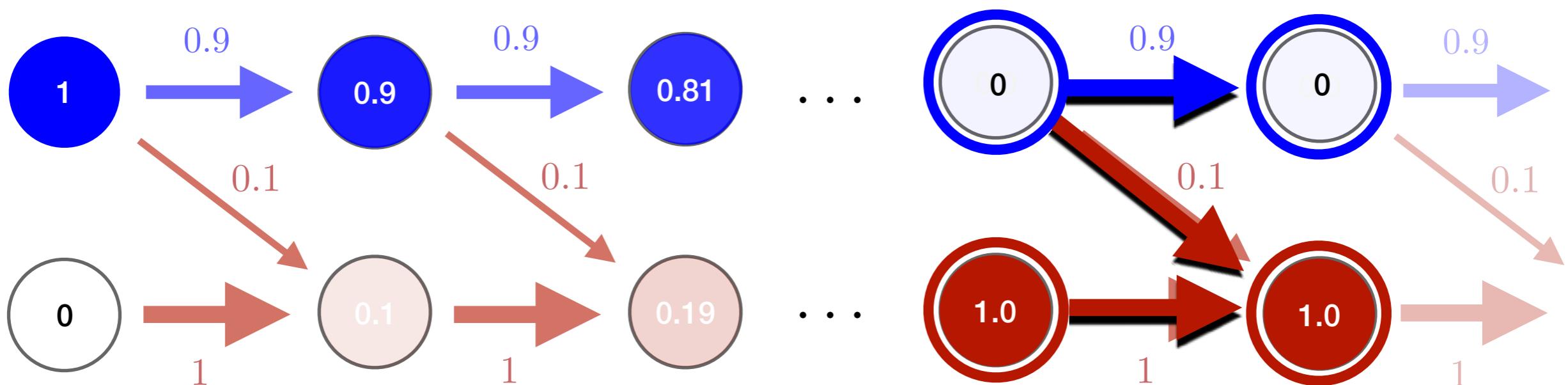
$t = 1$

$t = 2$

$t = 29$

$t = 30$

OK



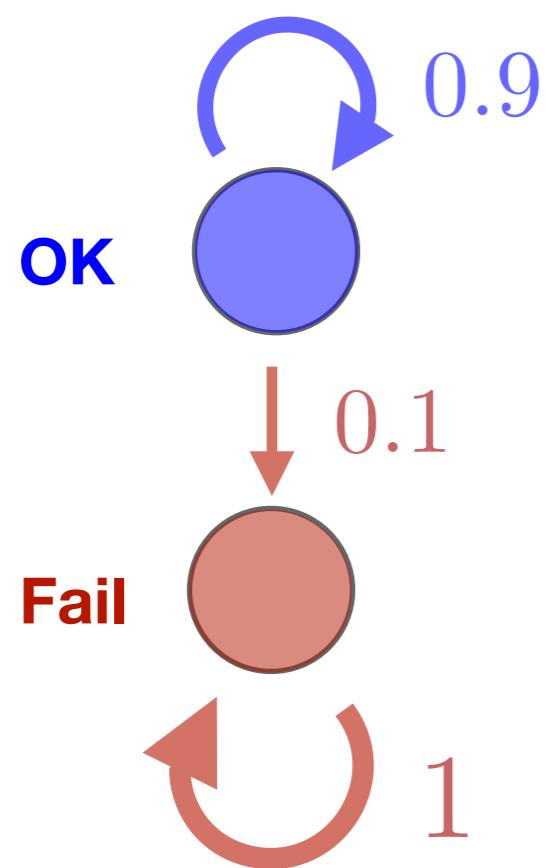
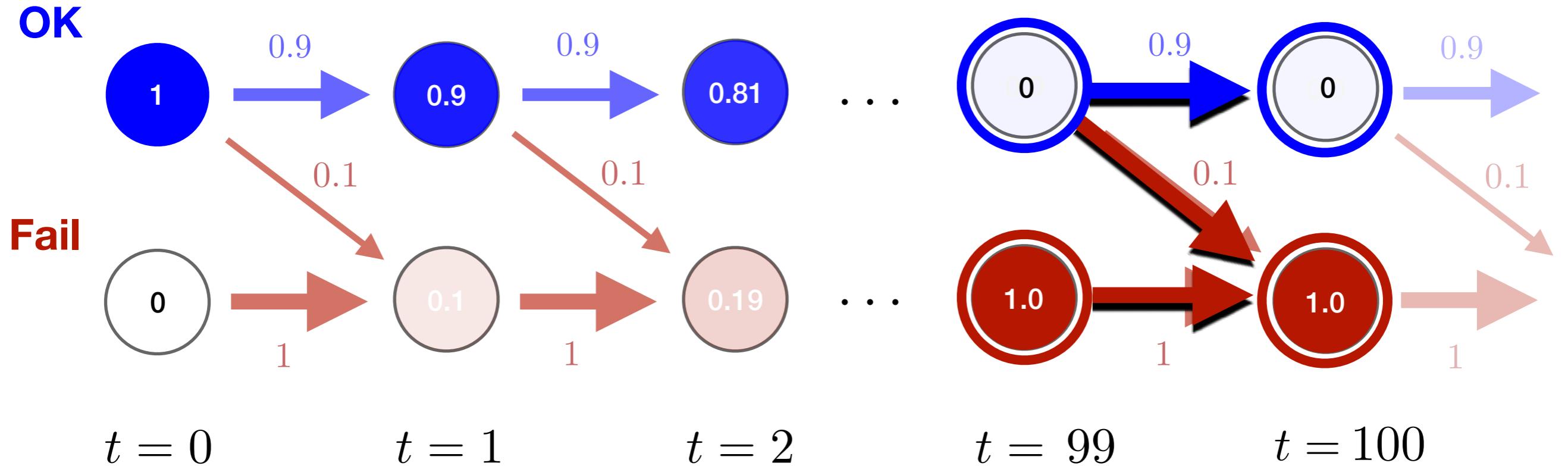
$t = 0$

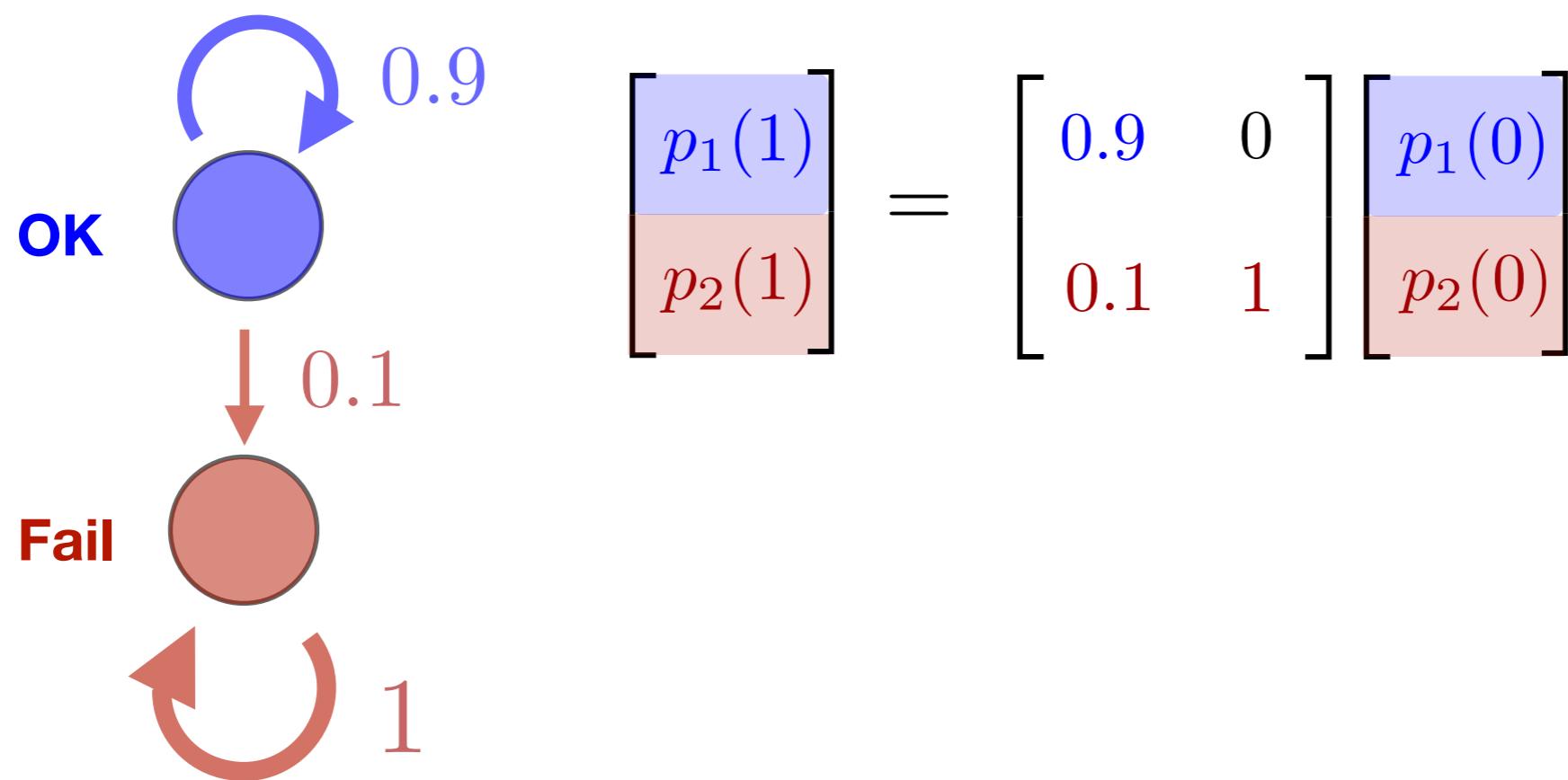
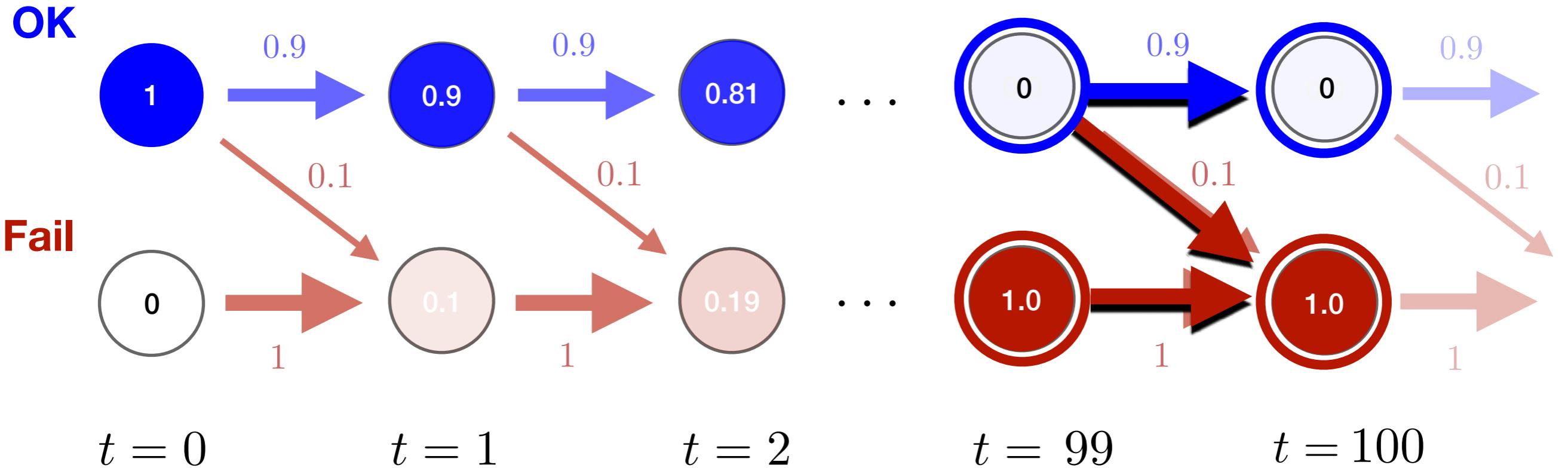
$t = 1$

$t = 2$

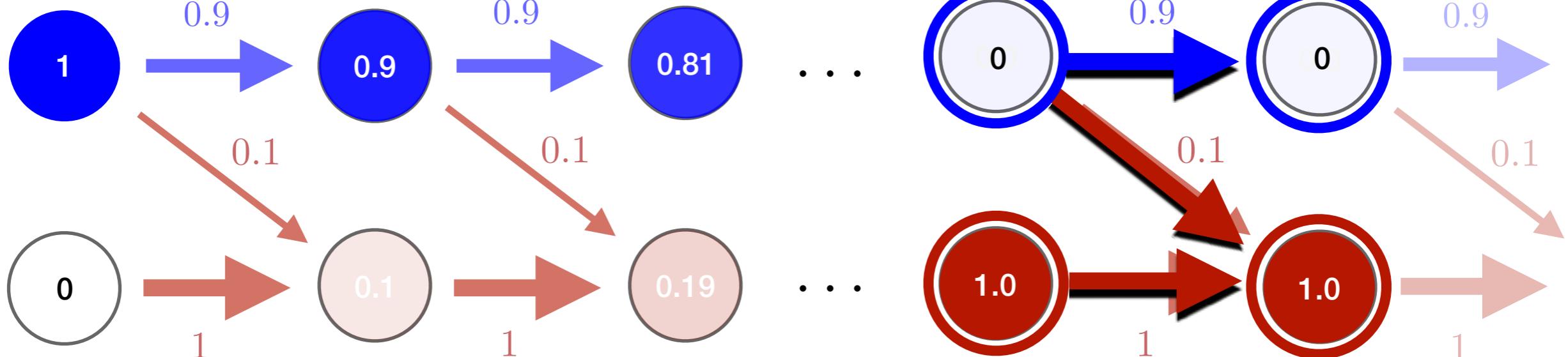
$t = 99$

$t = 100$





OK



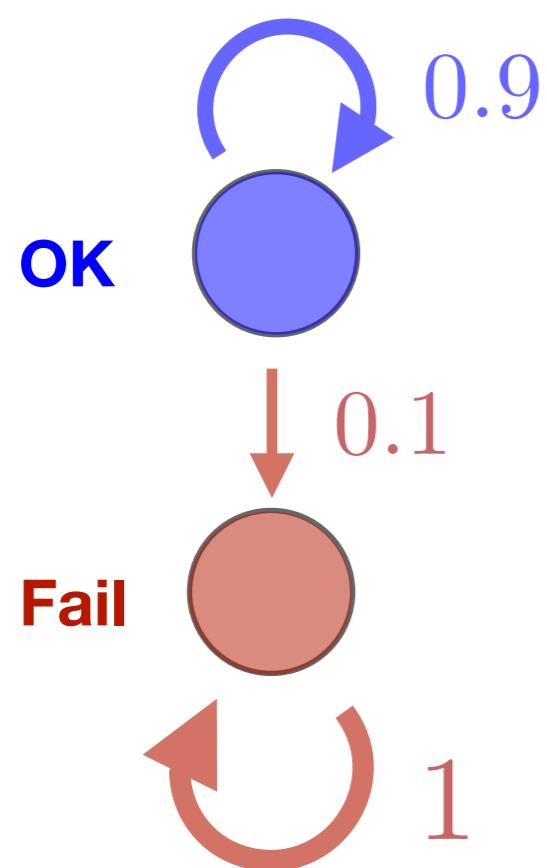
$t = 0$

$t = 1$

$t = 2$

$t = 99$

$t = 100$

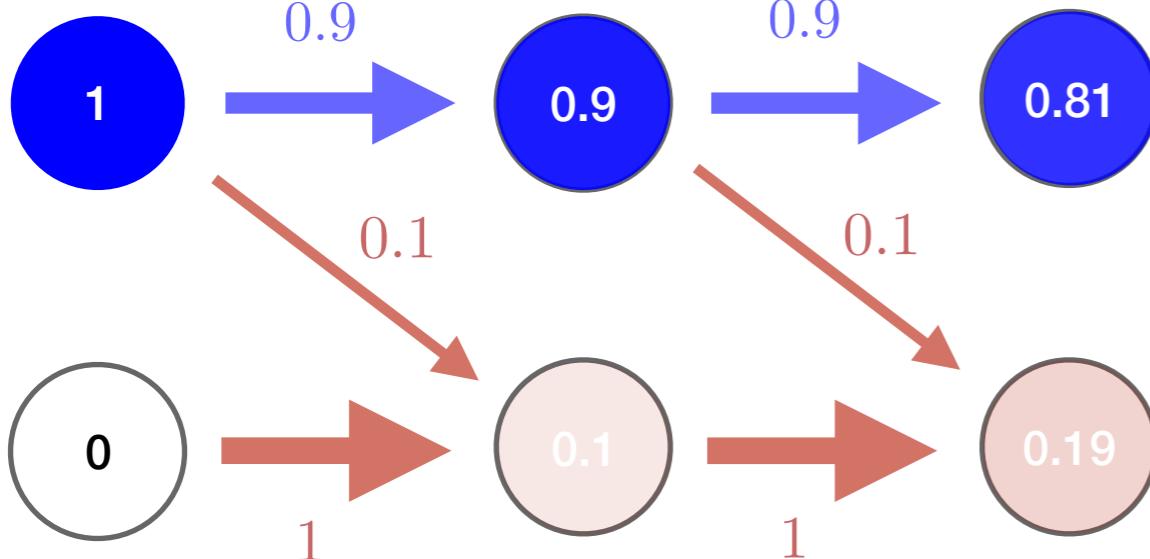


$$\begin{bmatrix} p_1(1) \\ p_2(1) \end{bmatrix} = \begin{bmatrix} 0.9 & 0 \\ 0.1 & 1 \end{bmatrix} \begin{bmatrix} p_1(0) \\ p_2(0) \end{bmatrix}$$

$$p_1(1) = 0.9 \times p_1(0) + 0 \times p_2(0)$$

$$p_2(1) = 0.1 \times p_1(0) + 1 \times p_2(0)$$

OK



$t = 0$

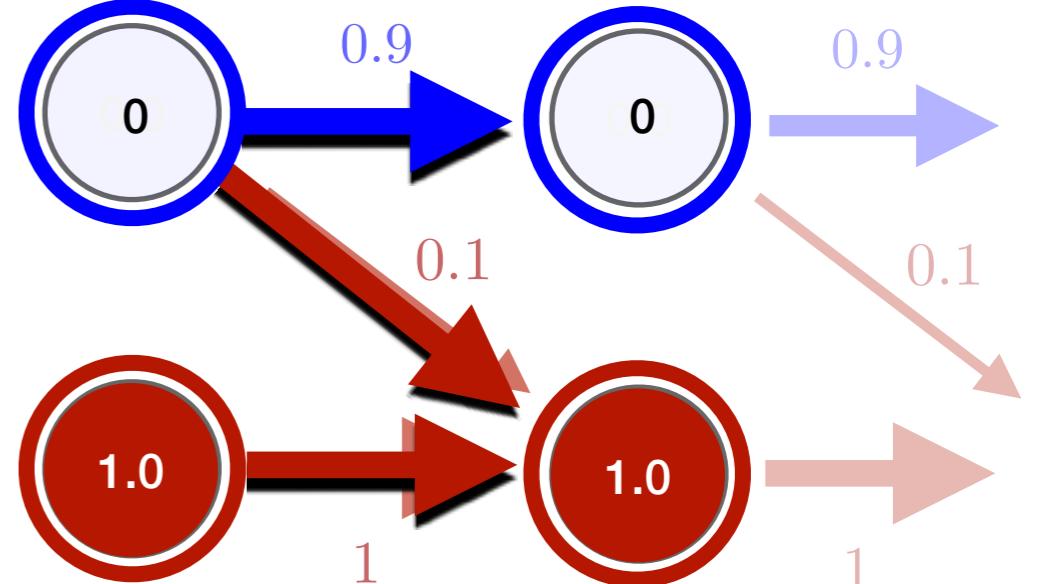
$t = 1$

$t = 2$

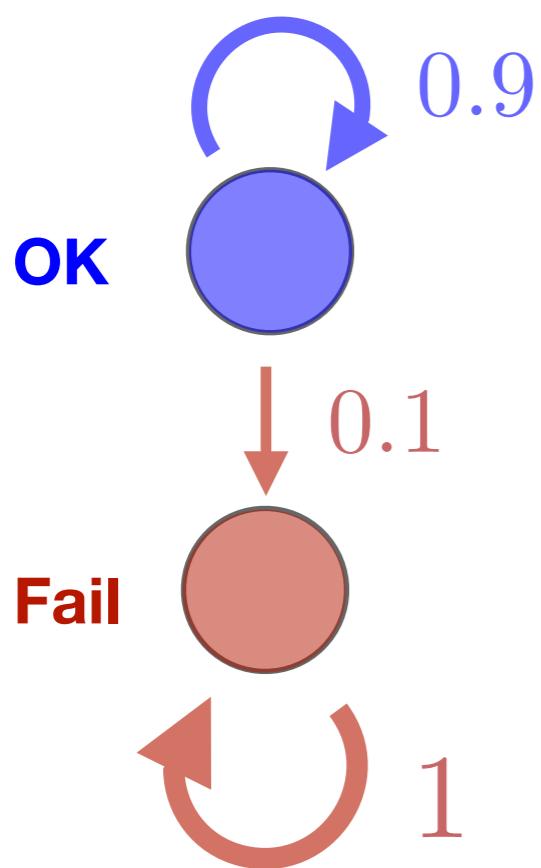
$t = 99$

$t = 100$

...



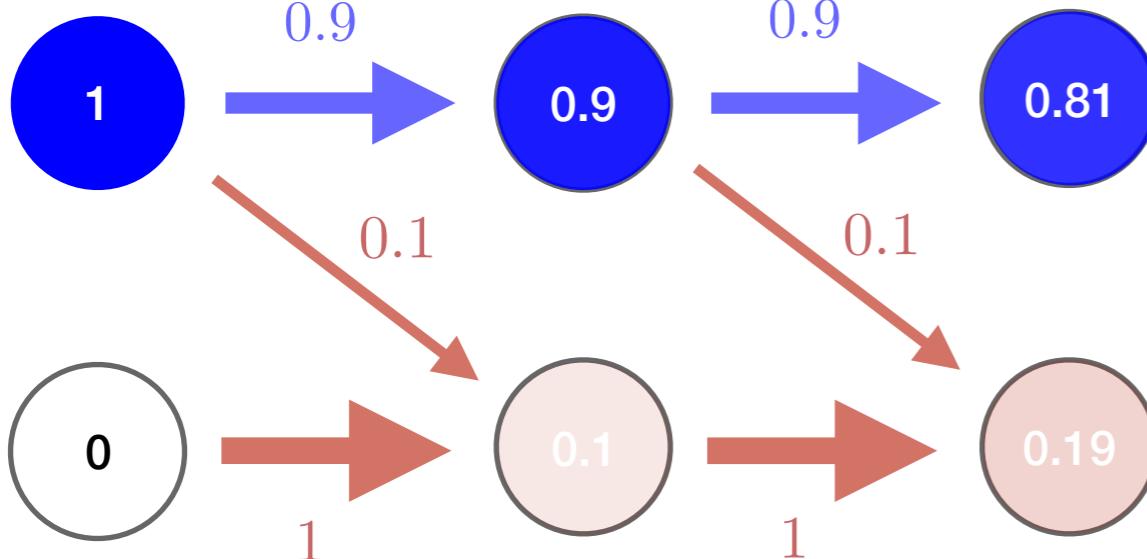
...



$$\begin{bmatrix} p_1(1) \\ p_2(1) \end{bmatrix} = \begin{bmatrix} 0.9 & 0 \\ 0.1 & 1 \end{bmatrix} \begin{bmatrix} p_1(0) \\ p_2(0) \end{bmatrix}$$

$$\begin{bmatrix} p_1(2) \\ p_2(2) \end{bmatrix} = \begin{bmatrix} 0.9 & 0 \\ 0.1 & 1 \end{bmatrix} \begin{bmatrix} p_1(1) \\ p_2(1) \end{bmatrix}$$

OK



$t = 0$

$t = 1$

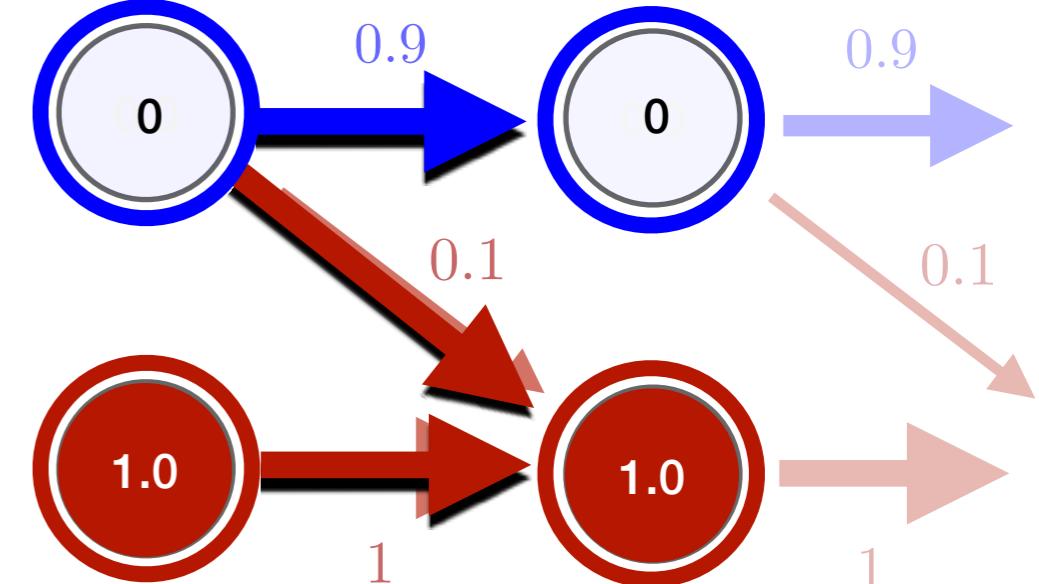
$t = 2$

$t = 99$

$t = 100$

Fail

...



0.1

0.1

1

1

0.1

1



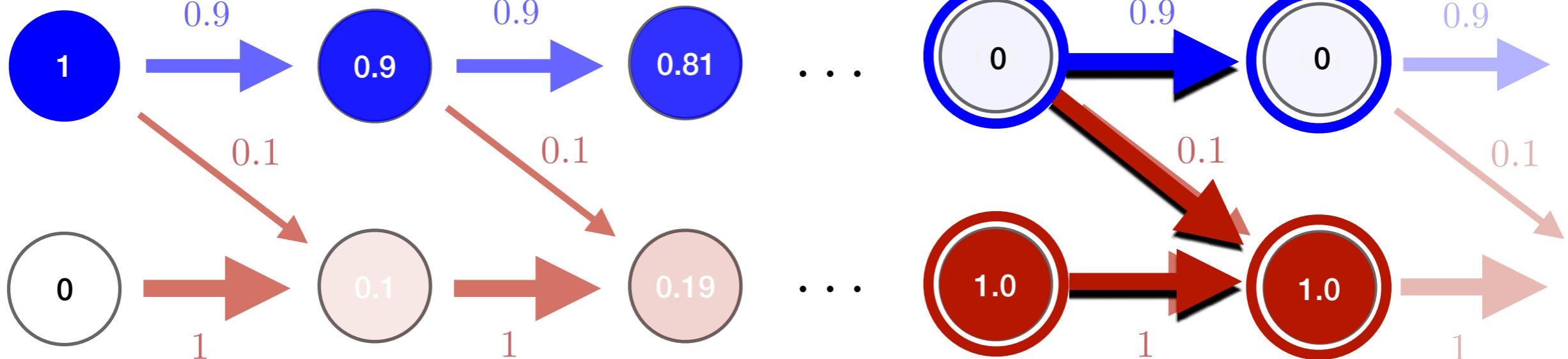
Fail



$$\begin{bmatrix} p_1(1) \\ p_2(1) \end{bmatrix} = \begin{bmatrix} 0.9 & 0 \\ 0.1 & 1 \end{bmatrix} \begin{bmatrix} p_1(0) \\ p_2(0) \end{bmatrix}$$

$$\begin{bmatrix} p_1(2) \\ p_2(2) \end{bmatrix} = \begin{bmatrix} 0.9 & 0 \\ 0.1 & 1 \end{bmatrix} \begin{bmatrix} 0.9 & 0 \\ 0.1 & 1 \end{bmatrix} \begin{bmatrix} p_1(0) \\ p_2(0) \end{bmatrix}$$

OK



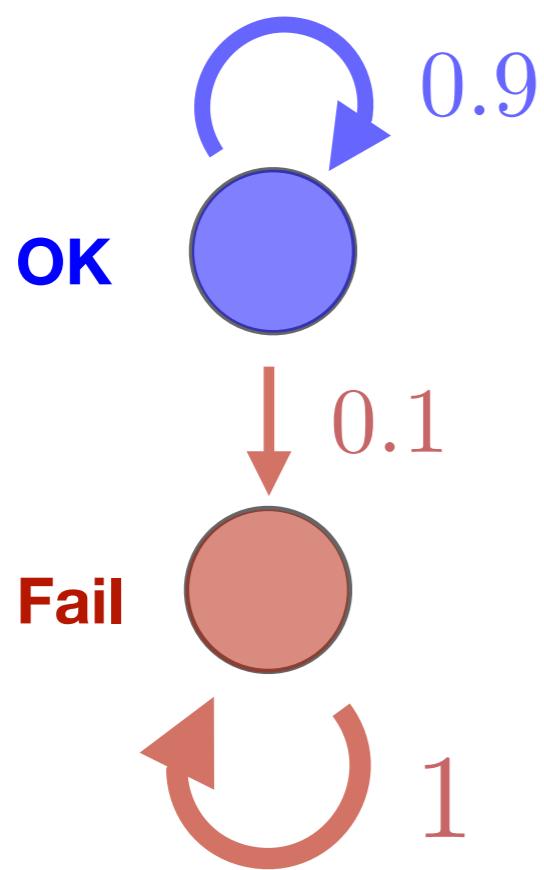
$t = 0$

$t = 1$

$t = 2$

$t = 99$

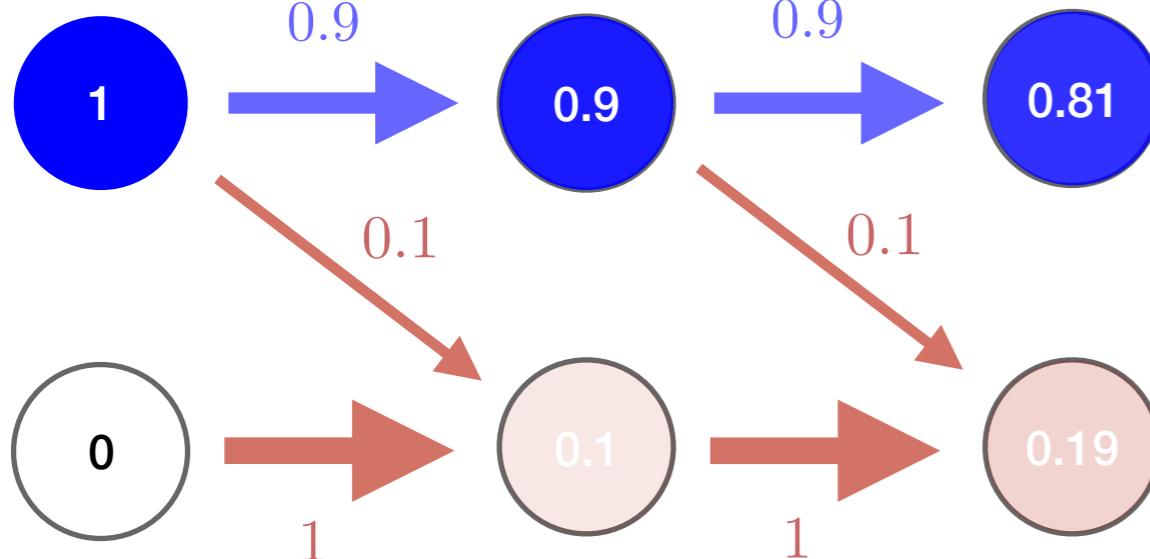
$t = 100$



$$\begin{bmatrix} p_1(1) \\ p_2(1) \end{bmatrix} = \begin{bmatrix} 0.9 & 0 \\ 0.1 & 1 \end{bmatrix} \begin{bmatrix} p_1(0) \\ p_2(0) \end{bmatrix}$$

$$\begin{bmatrix} p_1(2) \\ p_2(2) \end{bmatrix} = \begin{bmatrix} 0.9 & 0 \\ 0.1 & 1 \end{bmatrix}^2 \begin{bmatrix} p_1(0) \\ p_2(0) \end{bmatrix}$$

OK



$t = 0$

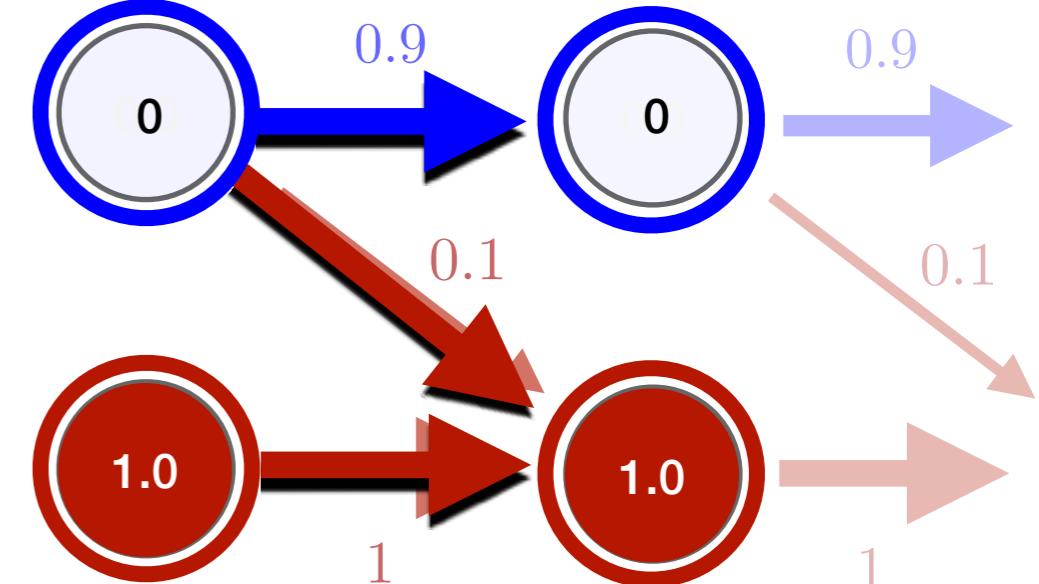
$t = 1$

$t = 2$

$t = 99$

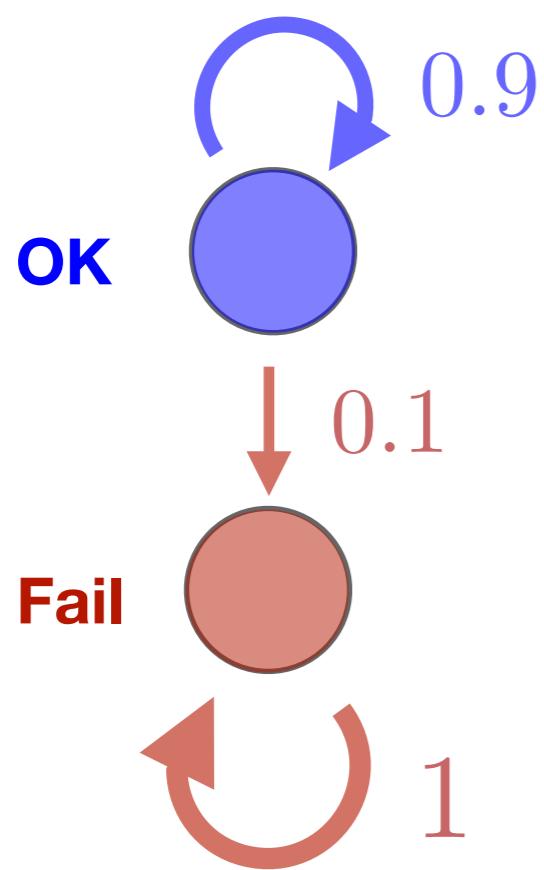
$t = 100$

...



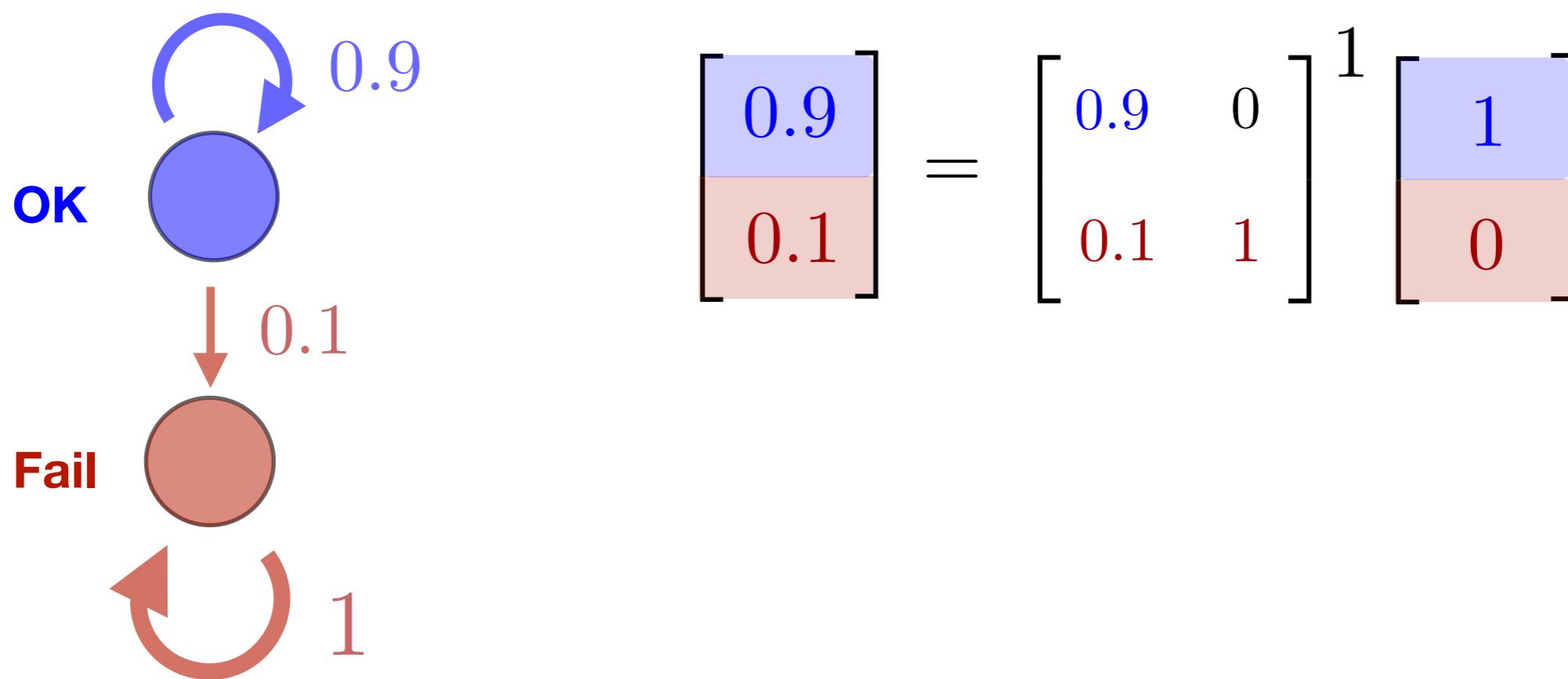
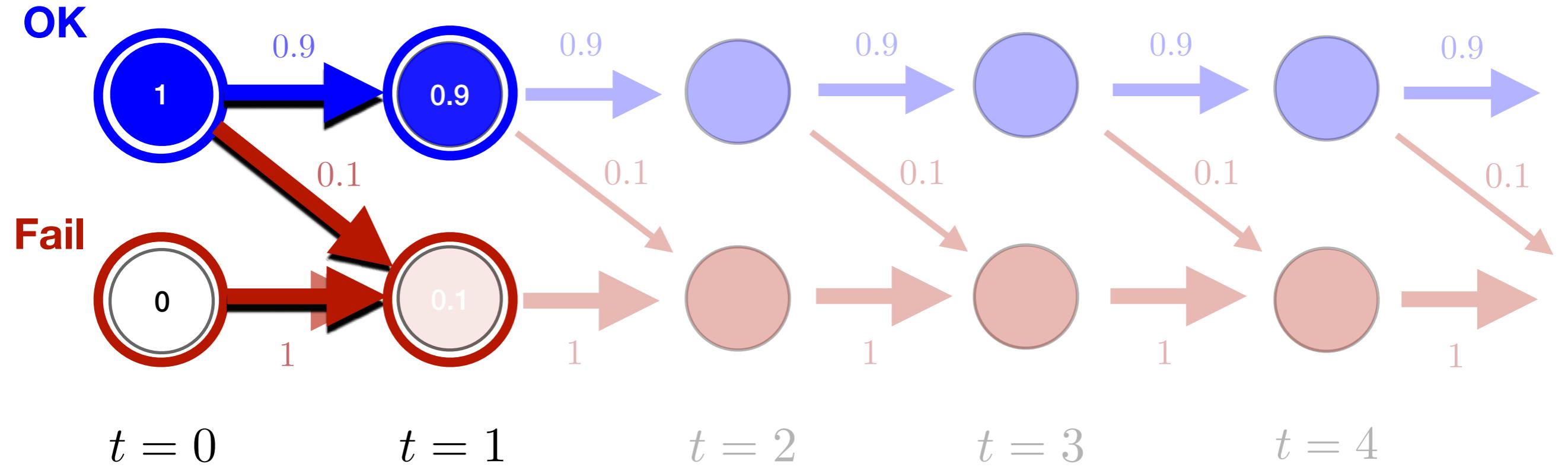
...

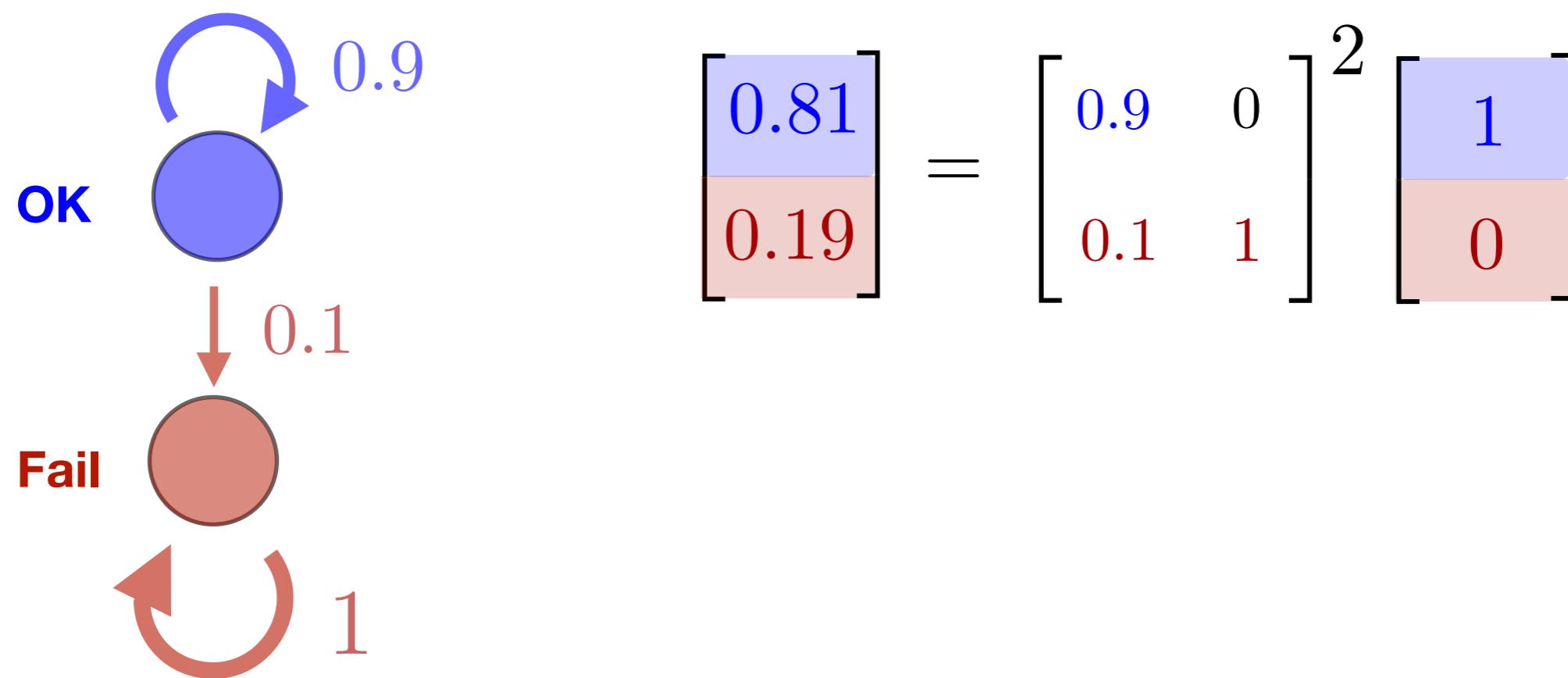
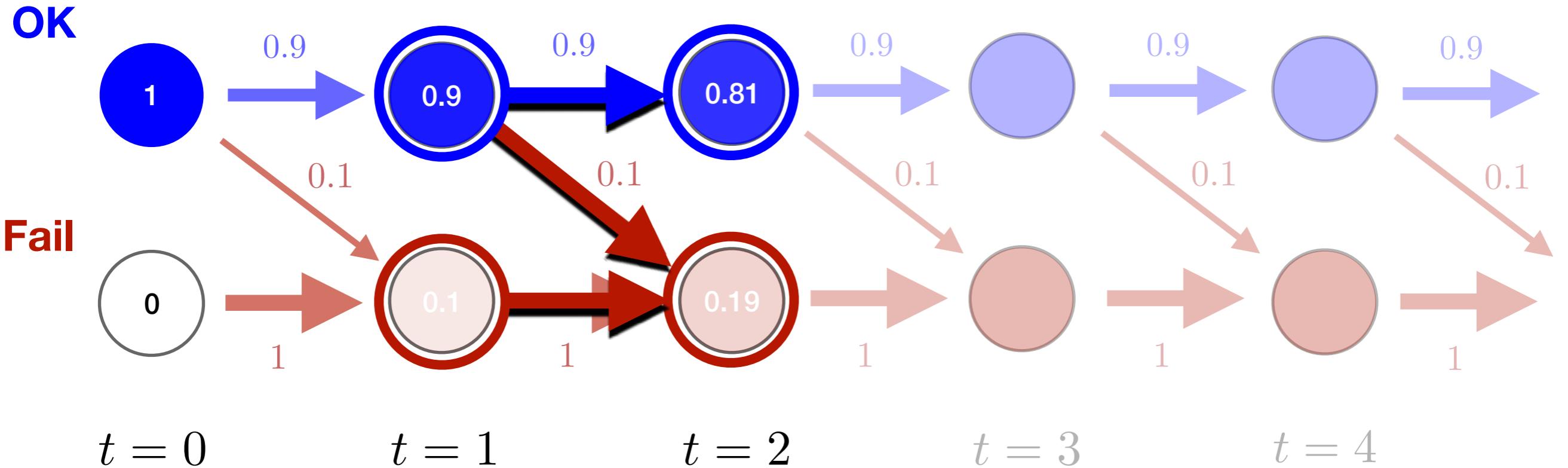
...

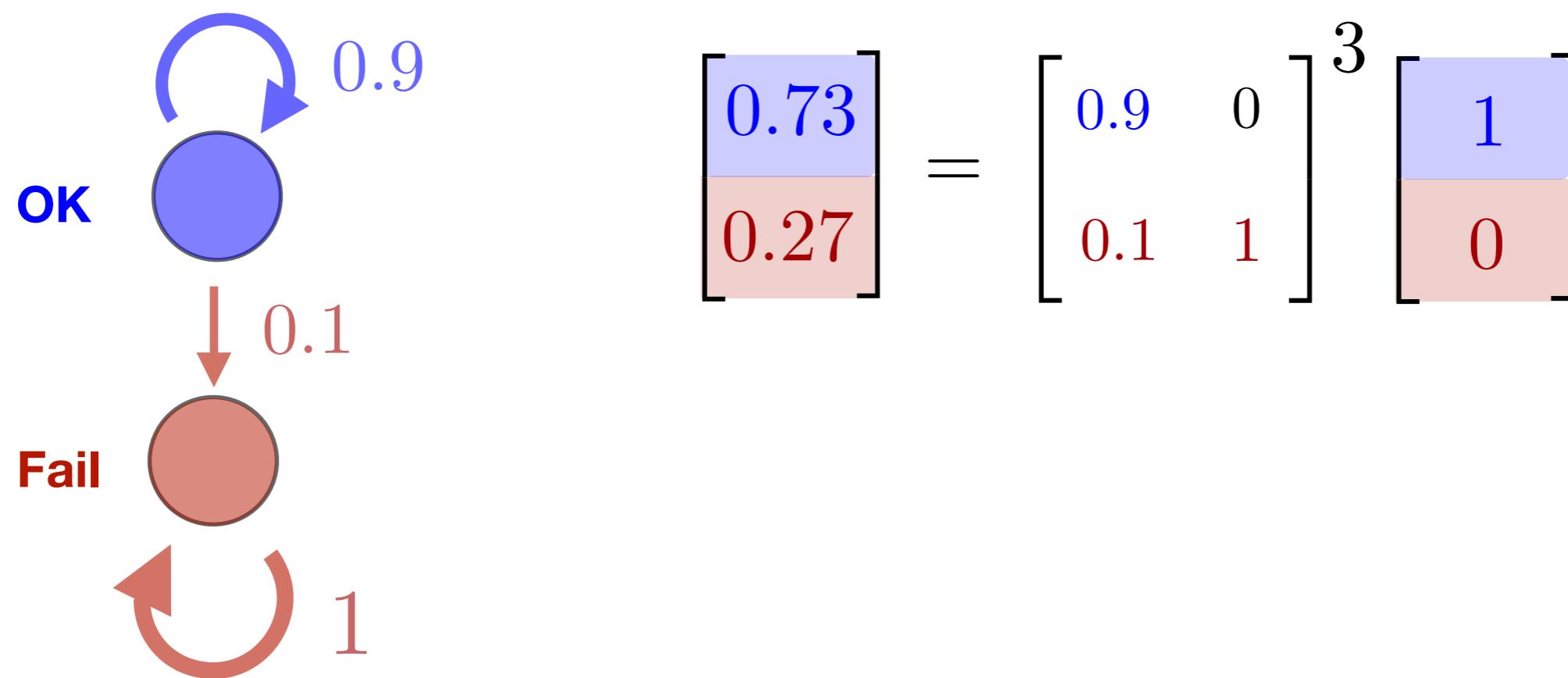
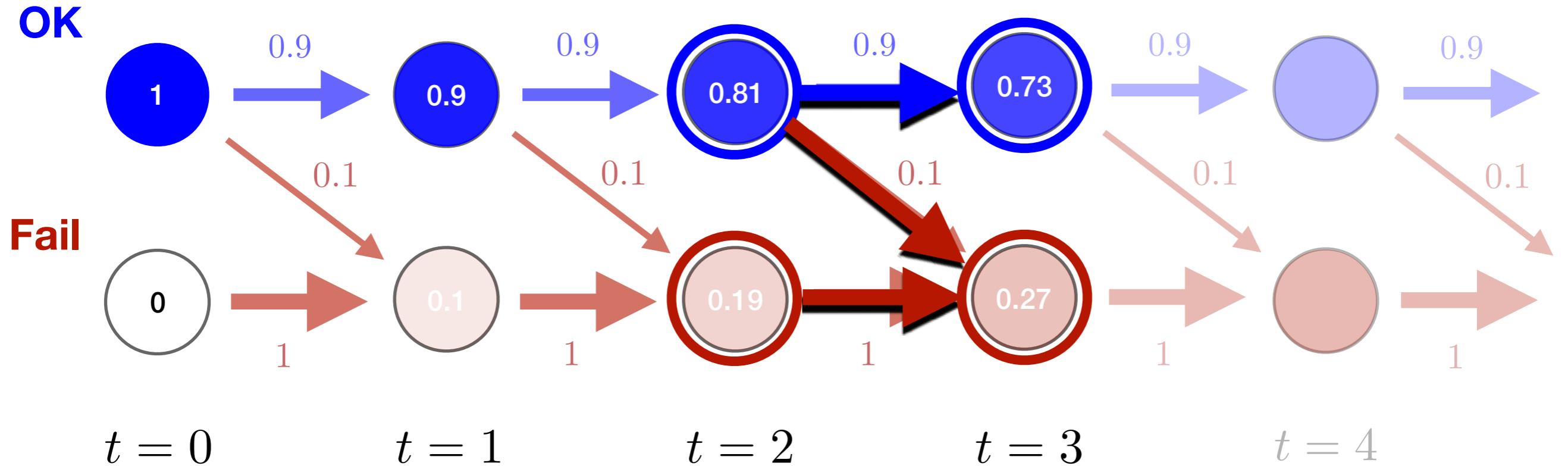


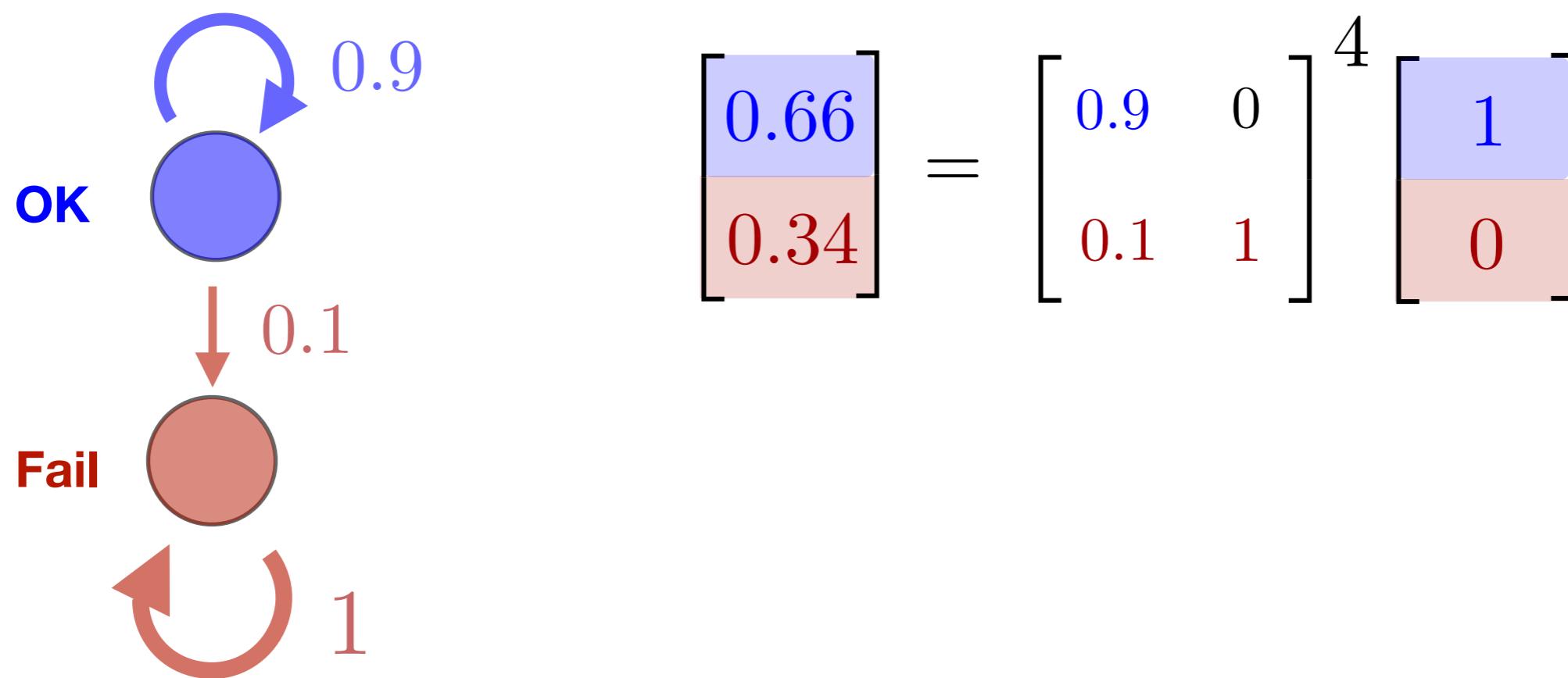
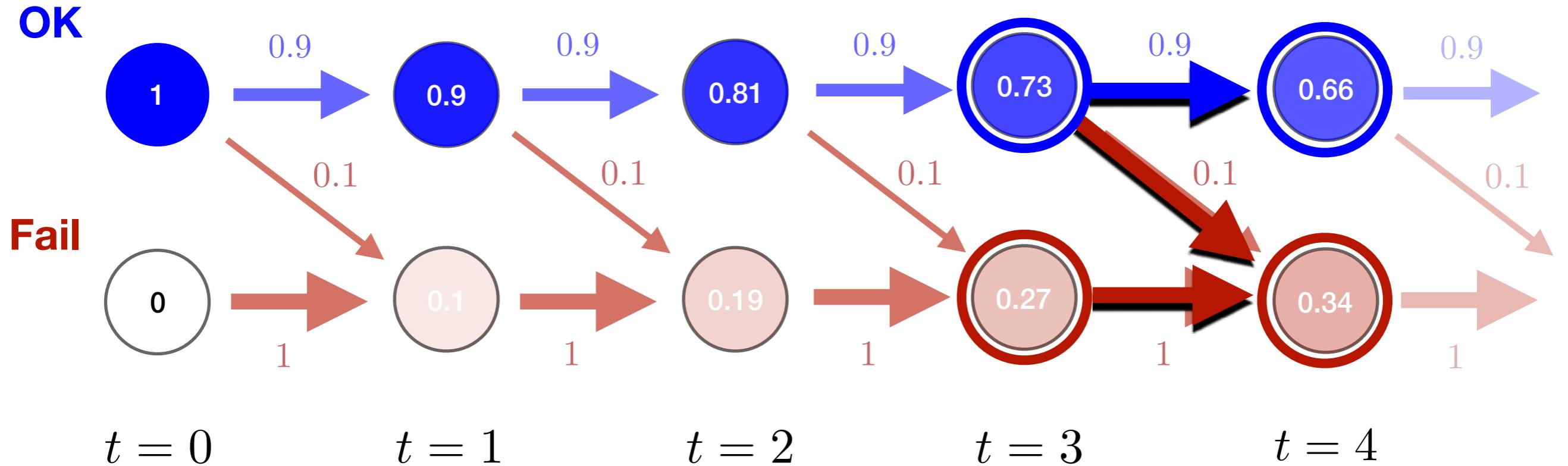
$$\begin{bmatrix} p_1(1) \\ p_2(1) \end{bmatrix} = \begin{bmatrix} 0.9 & 0 \\ 0.1 & 1 \end{bmatrix} \begin{bmatrix} p_1(0) \\ p_2(0) \end{bmatrix}$$

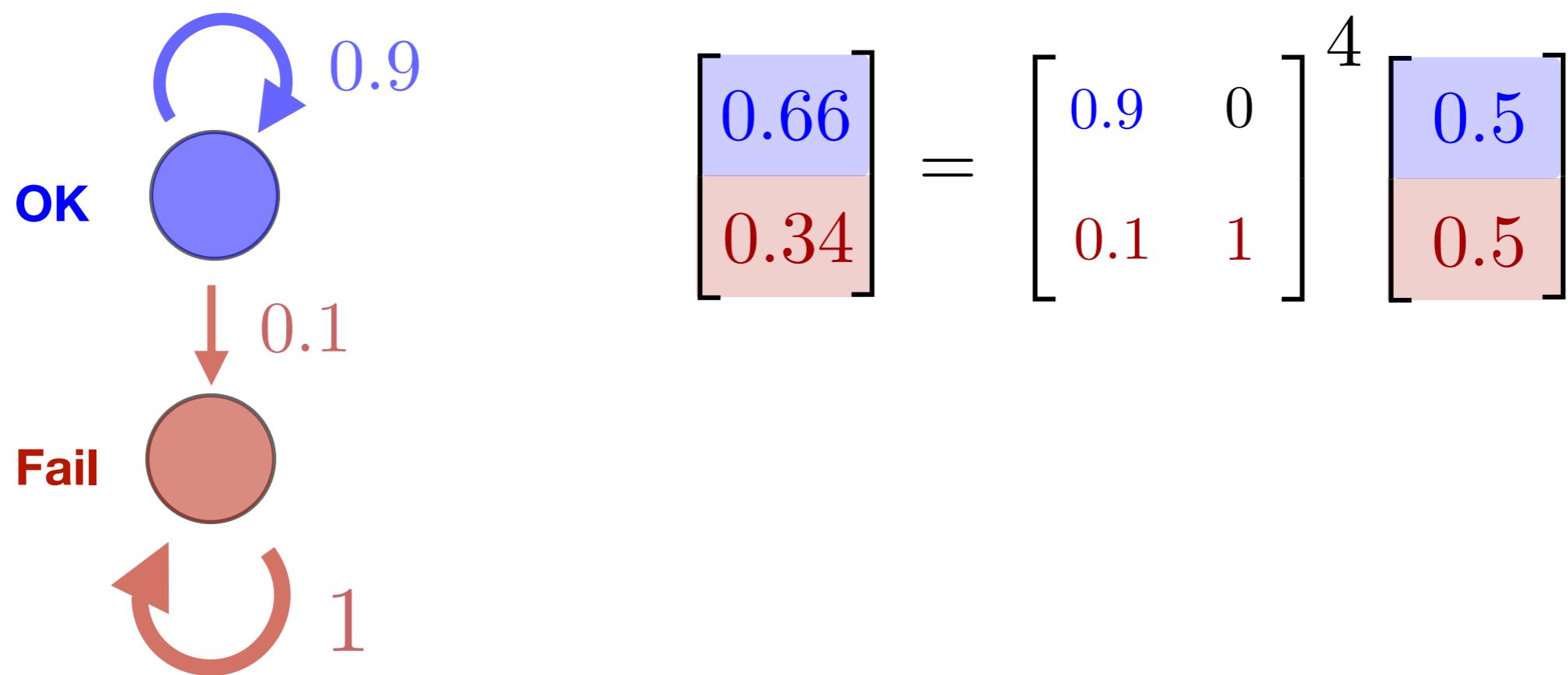
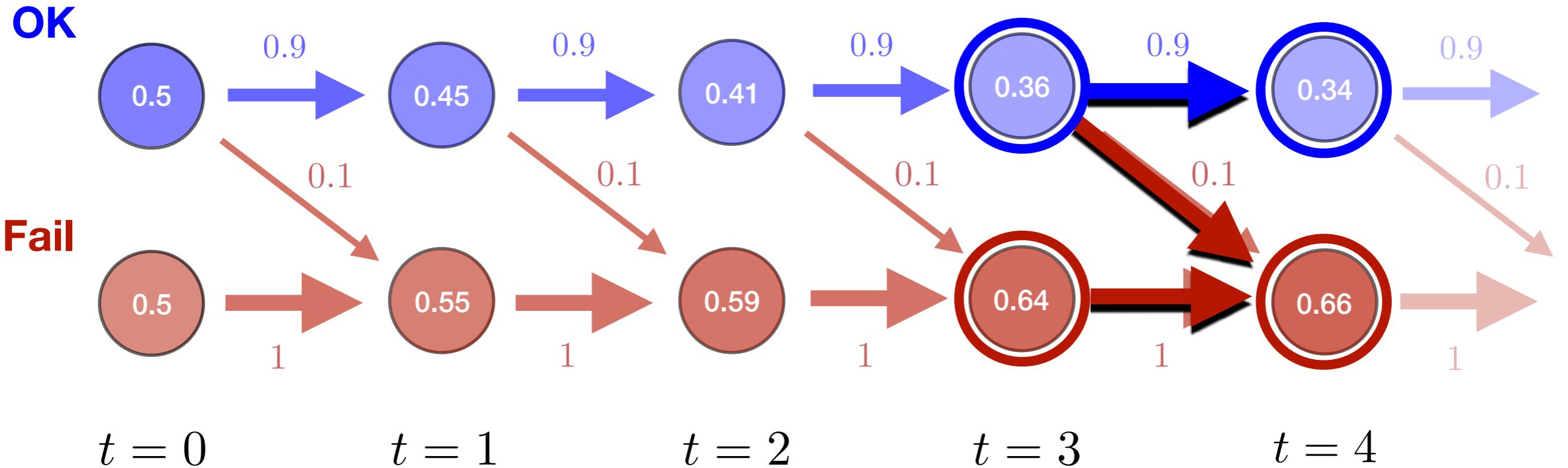
$$\begin{bmatrix} p_1(k) \\ p_2(k) \end{bmatrix} = \begin{bmatrix} 0.9 & 0 \\ 0.1 & 1 \end{bmatrix}^k \begin{bmatrix} p_1(0) \\ p_2(0) \end{bmatrix}$$

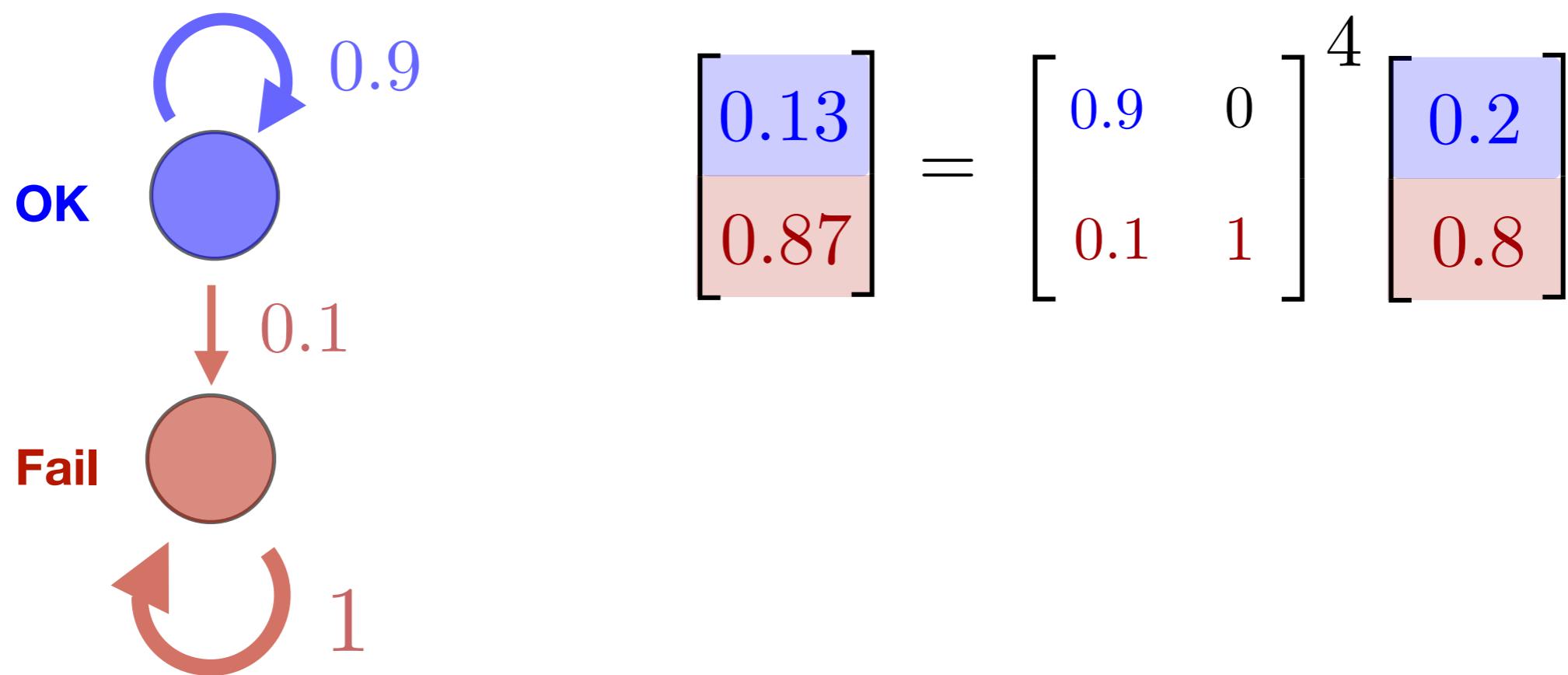
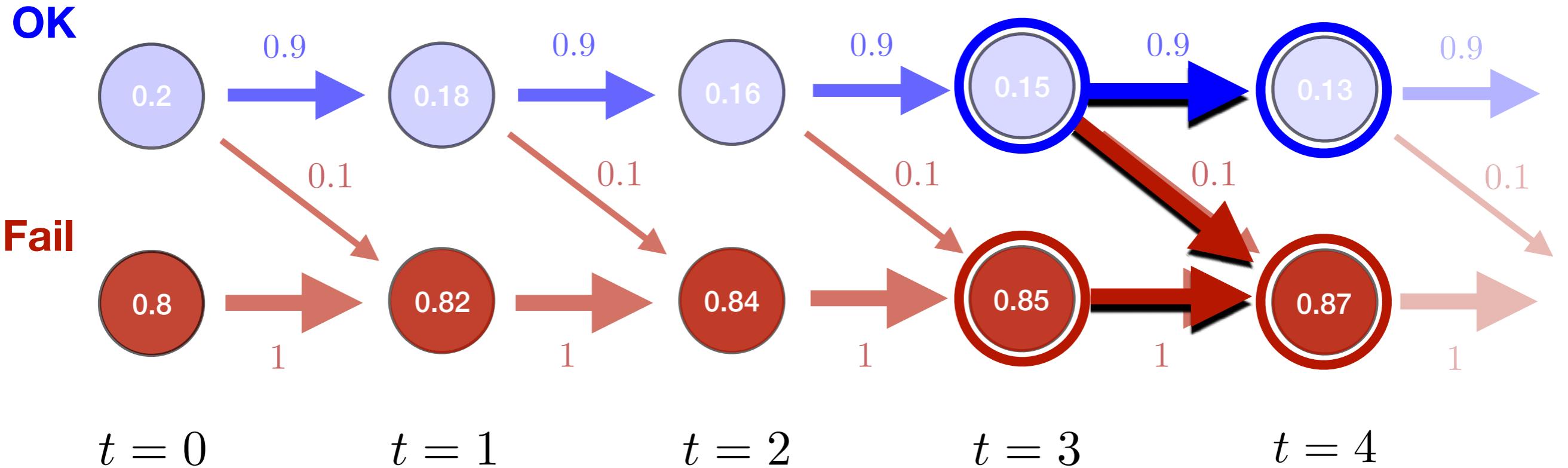




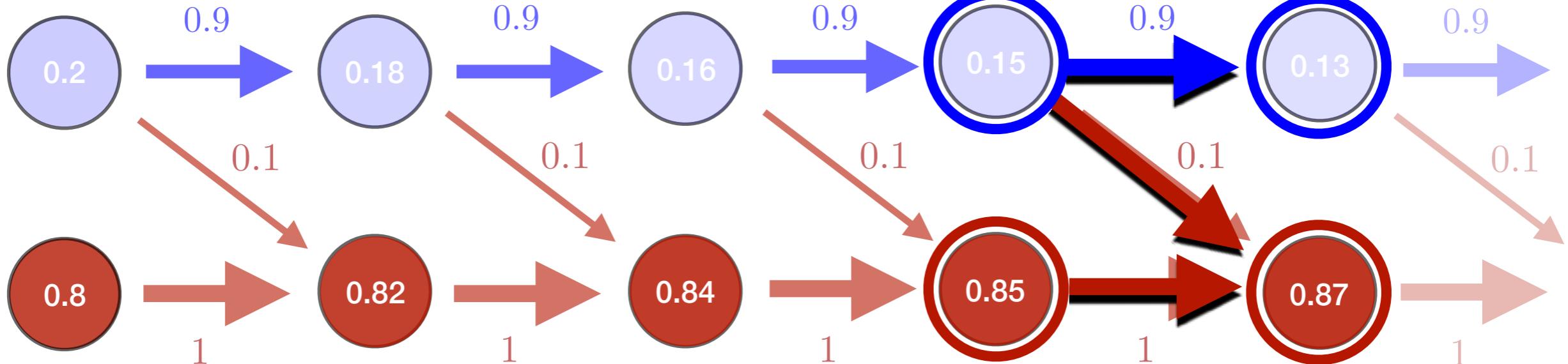




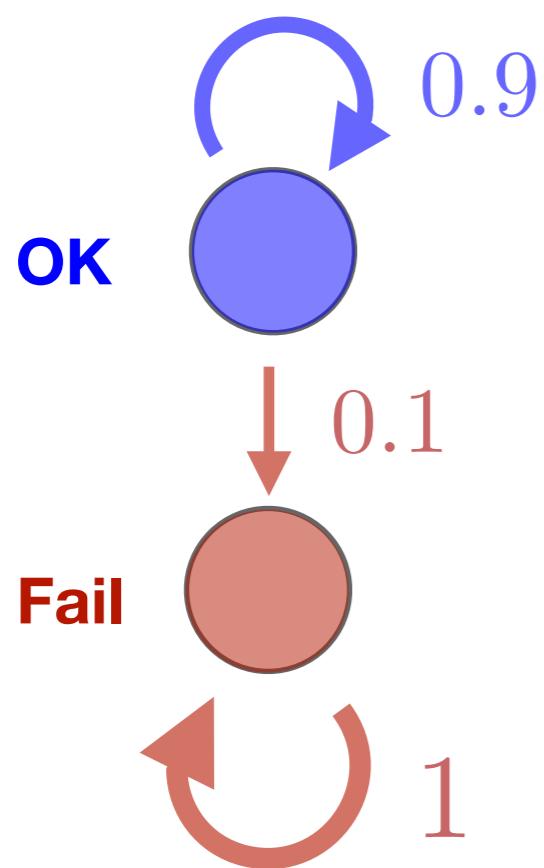




OK

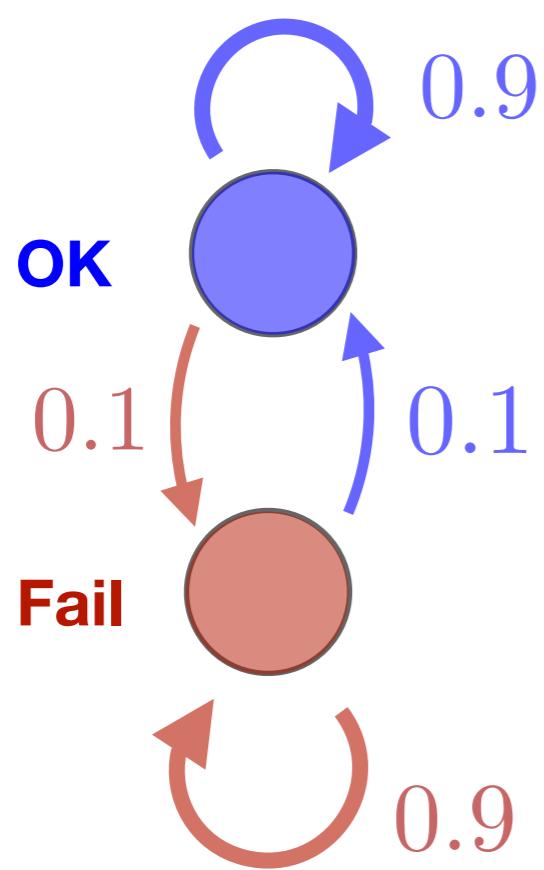
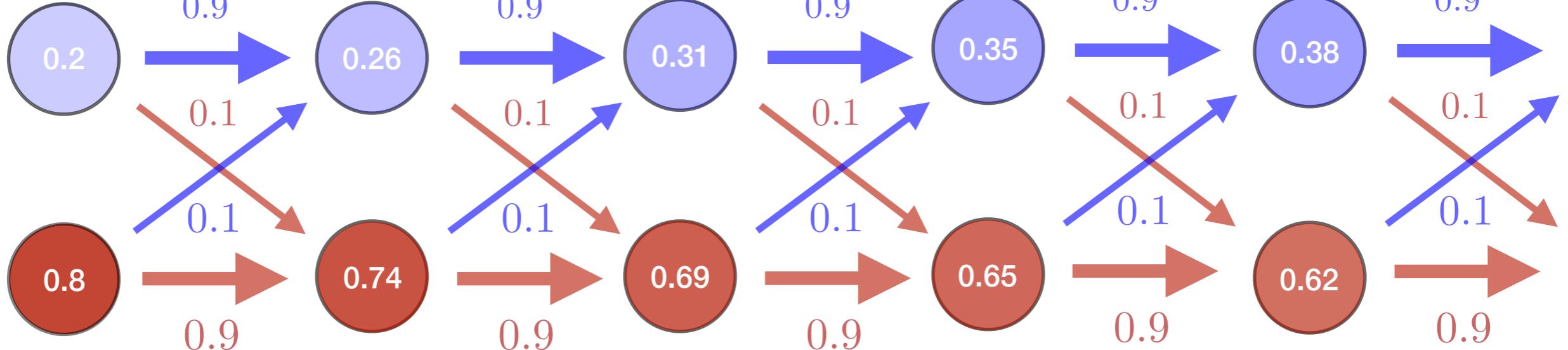


$t = 0 \quad t = 1 \quad t = 2 \quad t = 3 \quad t = 4$



$$\begin{bmatrix} p_1(k) \\ p_2(k) \end{bmatrix} = \begin{bmatrix} 0.9 & 0 \\ 0.1 & 1 \end{bmatrix}^k \begin{bmatrix} p_1(0) \\ p_2(0) \end{bmatrix}$$

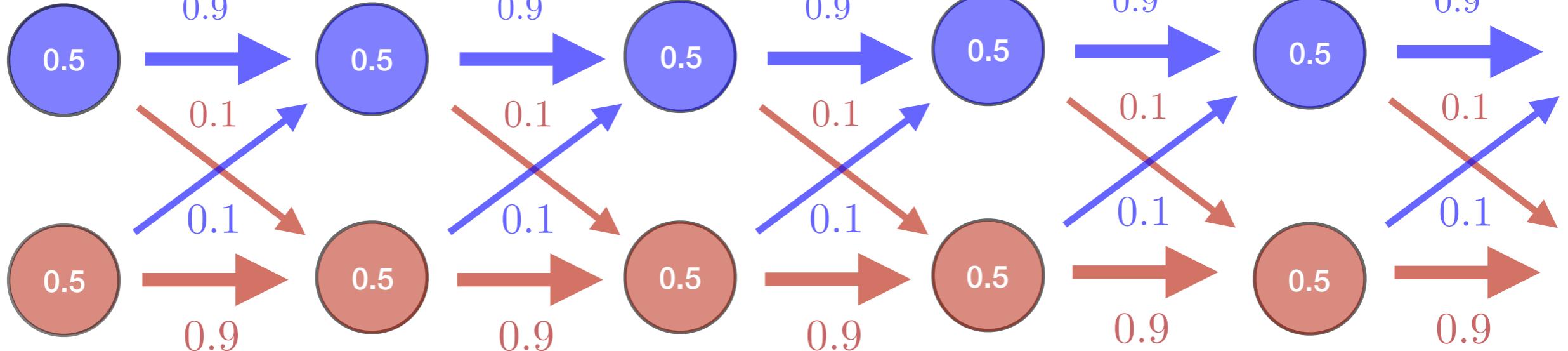
$$\lim_{k \rightarrow \infty} \begin{bmatrix} p_1(k) \\ p_2(k) \end{bmatrix} = \begin{bmatrix} p_1(\infty) \\ p_2(\infty) \end{bmatrix} = \begin{bmatrix} 0.0 \\ 1.0 \end{bmatrix}$$

OK

$$\begin{bmatrix} p_1(k) \\ p_2(k) \end{bmatrix} = \begin{bmatrix} 0.9 & 0.1 \\ 0.1 & 0.9 \end{bmatrix}^k \begin{bmatrix} p_1(0) \\ p_2(0) \end{bmatrix}$$

$$\lim_{k \rightarrow \infty} \begin{bmatrix} p_1(k) \\ p_2(k) \end{bmatrix} = \begin{bmatrix} p_1(\infty) \\ p_2(\infty) \end{bmatrix} = \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix}$$

OK



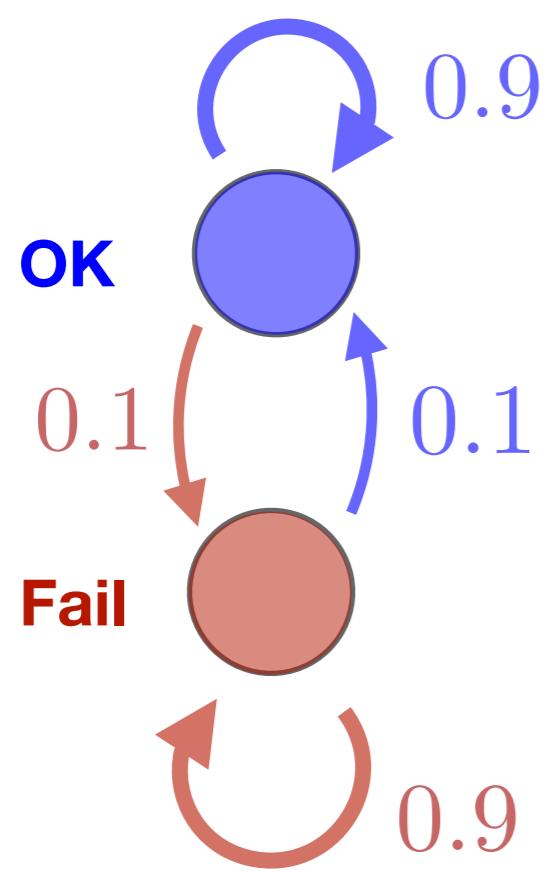
$t = 0$

$t = 1$

$t = 2$

$t = 3$

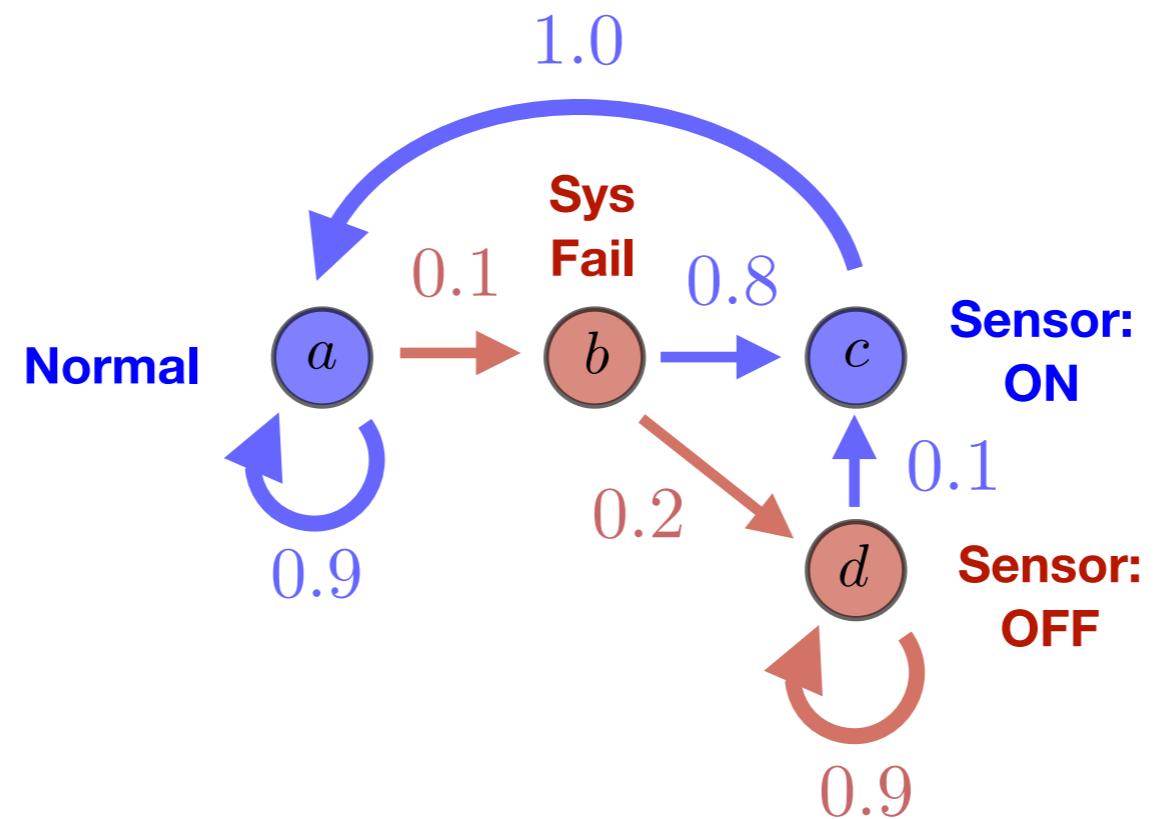
$t = 4$

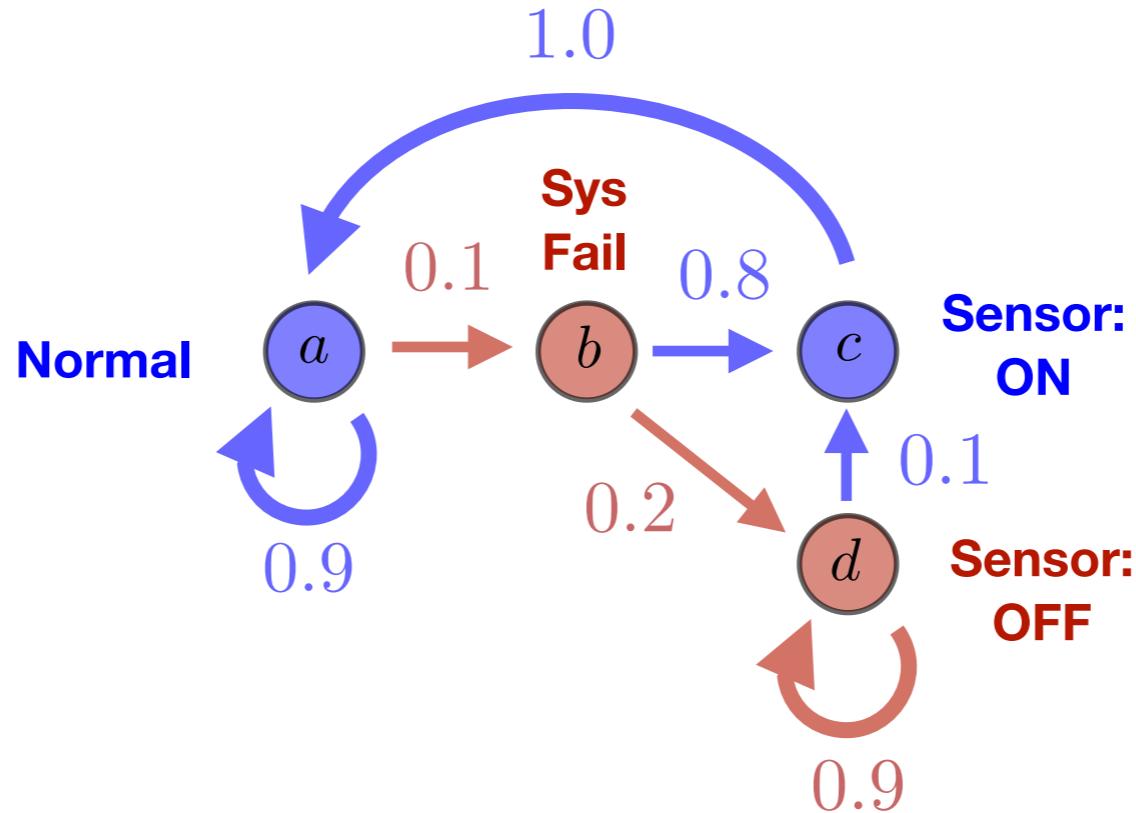


$$\begin{bmatrix} p_1(k) \\ p_2(k) \end{bmatrix} = \begin{bmatrix} 0.9 & 0.1 \\ 0.1 & 0.9 \end{bmatrix}^k \begin{bmatrix} p_1(0) \\ p_2(0) \end{bmatrix}$$

$$\lim_{k \rightarrow \infty} \begin{bmatrix} p_1(k) \\ p_2(k) \end{bmatrix} = \begin{bmatrix} p_1(\infty) \\ p_2(\infty) \end{bmatrix} = \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix}$$

examples



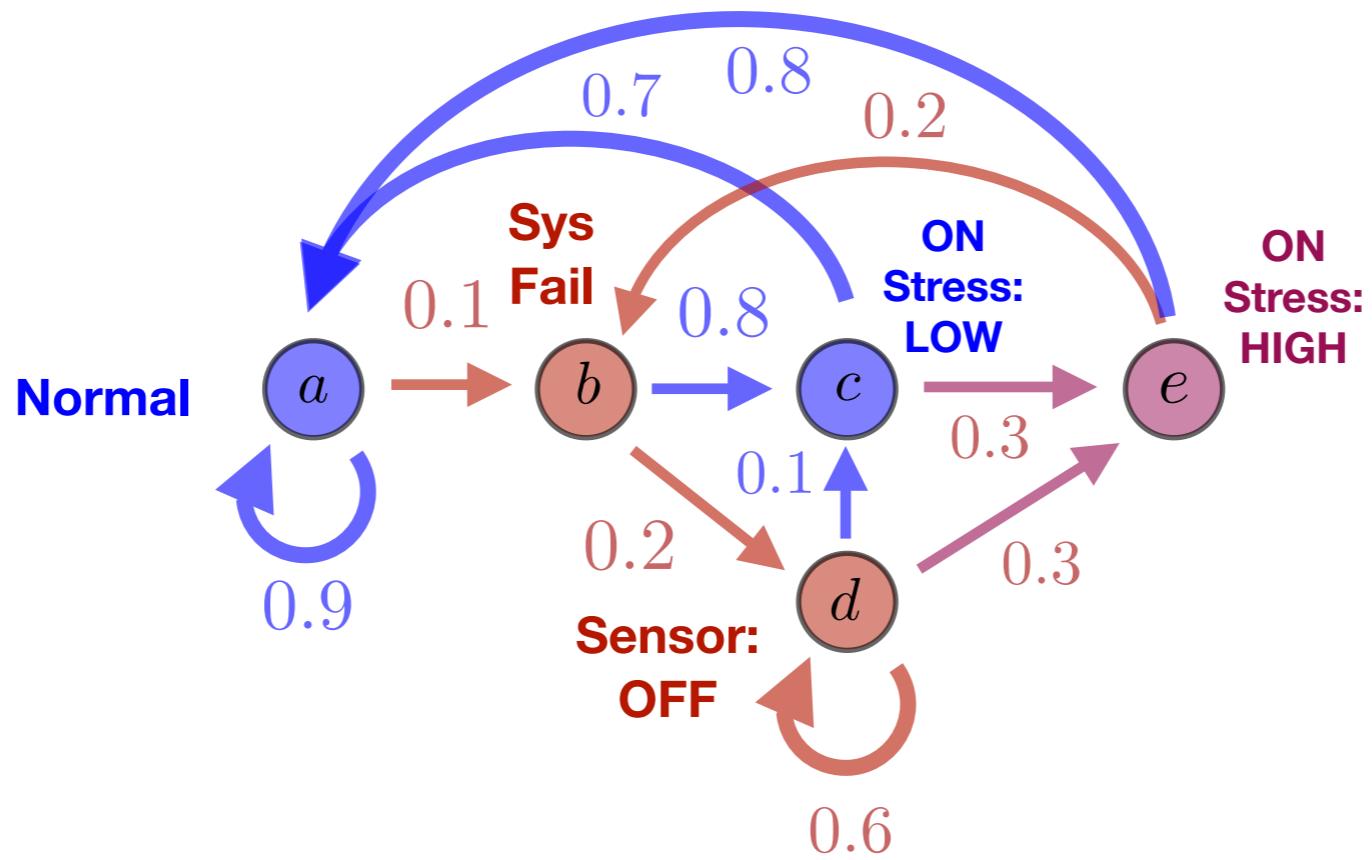


Update Equation:

$$\begin{bmatrix} p_a(k) \\ p_b(k) \\ p_c(k) \\ p_d(k) \end{bmatrix} = \begin{bmatrix} 0.9 & 0 & 1.0 & 0 \\ 0.1 & 0 & 0 & 0 \\ 0 & 0.8 & 0 & 0.1 \\ 0 & 0.2 & 0 & 0.9 \end{bmatrix}^k \begin{bmatrix} p_a(0) \\ p_b(0) \\ p_c(0) \\ p_d(0) \end{bmatrix}$$

Steady State Dist.

$$\begin{bmatrix} p_a(\infty) \\ p_b(\infty) \\ p_c(\infty) \\ p_d(\infty) \end{bmatrix} = \begin{bmatrix} 0.71 \\ 0.07 \\ 0.07 \\ 0.14 \end{bmatrix}$$



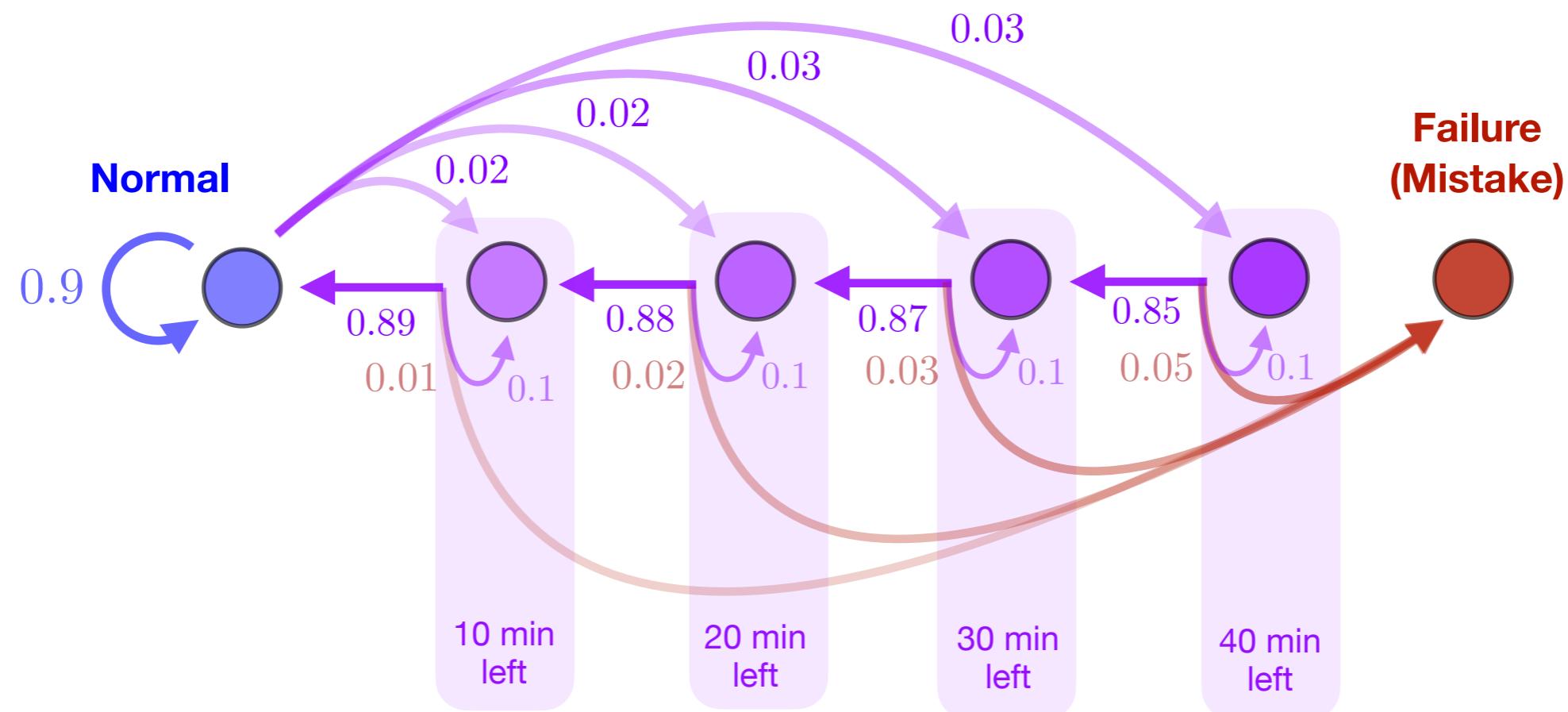
Update Equation:

Steady State Dist.

$$\begin{bmatrix} p_a(k) \\ p_b(k) \\ p_c(k) \\ p_d(k) \\ p_e(k) \end{bmatrix} = \begin{bmatrix} 0.9 & 0 & 0.7 & 0 & 0.8 \\ 0.1 & 0 & 0 & 0 & 0.2 \\ 0 & 0.8 & 0 & 0.1 & 0 \\ 0 & 0.2 & 0 & 0.6 & 0 \\ 0 & 0 & 0.3 & 0.3 & 0 \end{bmatrix}^k \begin{bmatrix} p_a(0) \\ p_b(0) \\ p_c(0) \\ p_d(0) \\ p_e(0) \end{bmatrix} = \begin{bmatrix} p_a(\infty) \\ p_b(\infty) \\ p_c(\infty) \\ p_d(\infty) \\ p_e(\infty) \end{bmatrix} = \begin{bmatrix} 0.77 \\ 0.08 \\ 0.07 \\ 0.04 \\ 0.03 \end{bmatrix}$$

Markov Chain Models

1. Determining model structure for specific scenarios.
2. Determine transition probabilities from HF considerations.
3. Compute failure probabilities.

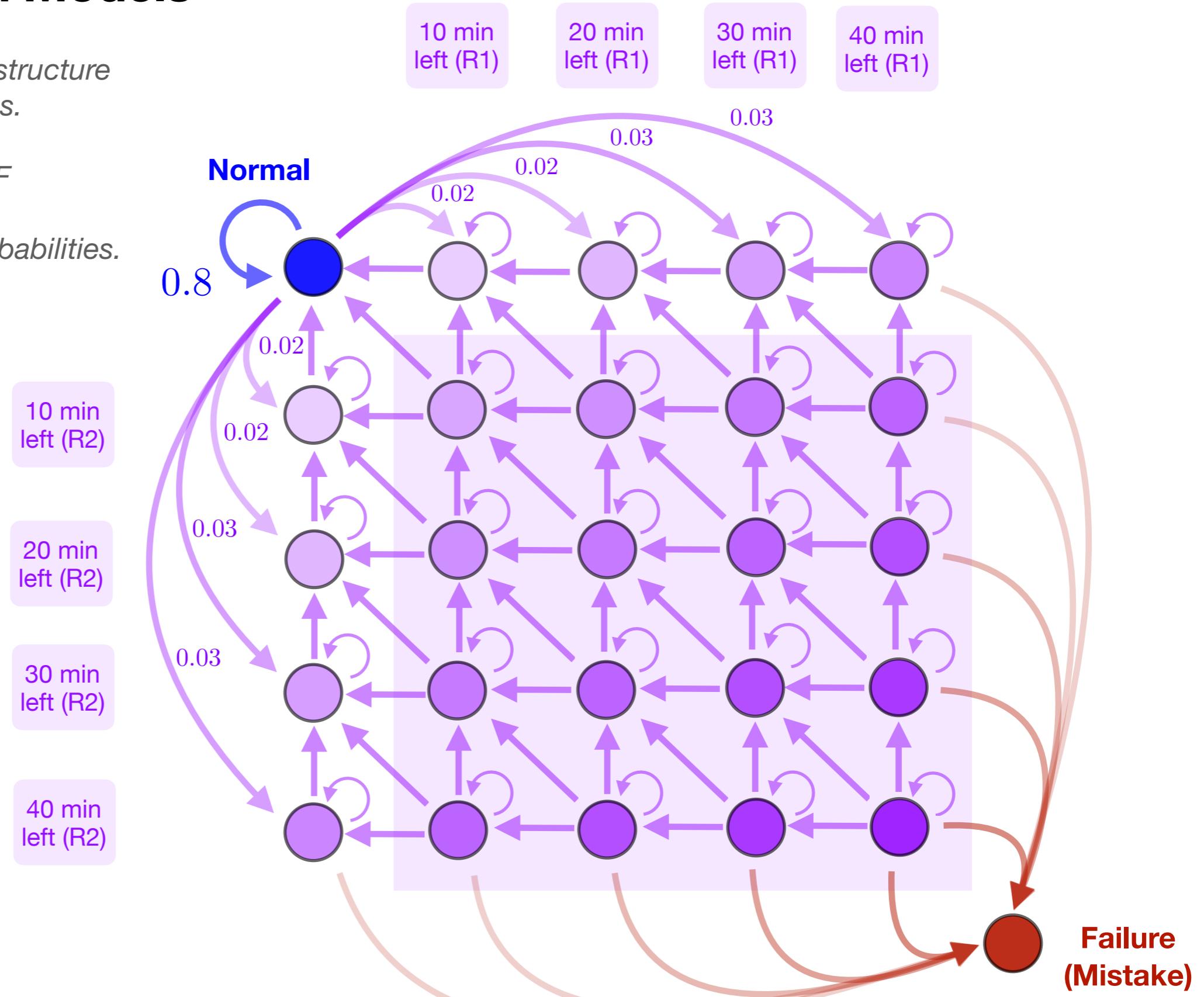


Markov Chain Models

1. Determining model structure for specific scenarios.
2. Determine transition probabilities from HF considerations.
3. Compute failure probabilities.

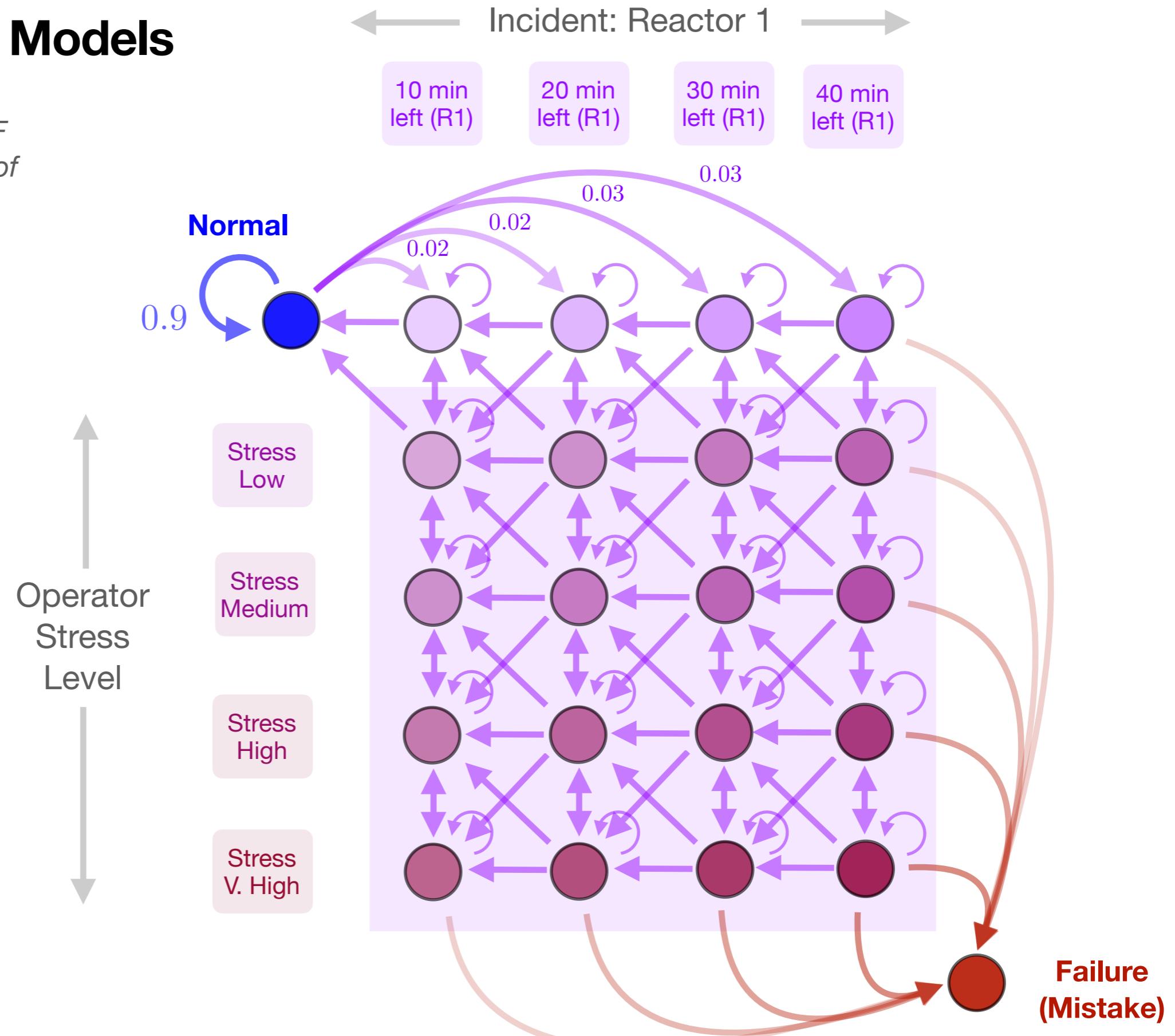
↑
Incident: Reactor 2
↓

← Incident: Reactor 1 →



Markov Chain Models

4. Modeling different HF components as part of the state space.

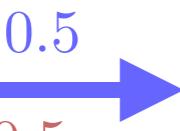
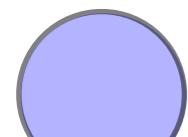


Markov Decision Process (MDP): Examples

Stochastic Processes

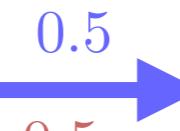
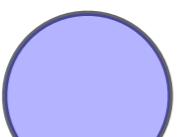
Major sources:

OK

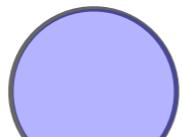


0.5

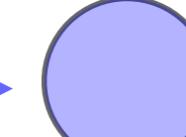
0.5



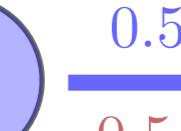
0.5



0.5

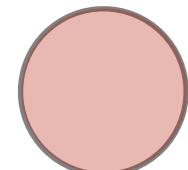


0.5



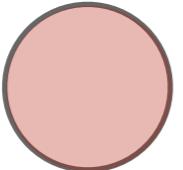
0.5

Fail



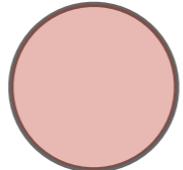
1

1



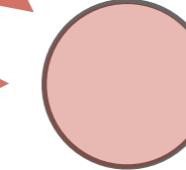
1

1



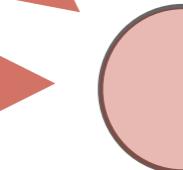
1

1



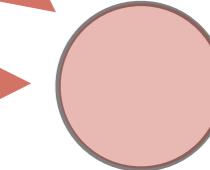
1

1



1

1



$t = 0$

$t = 1$

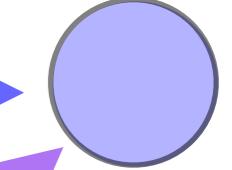
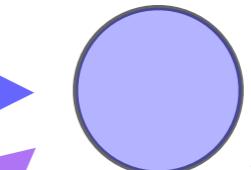
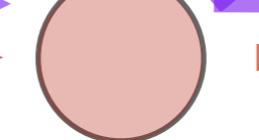
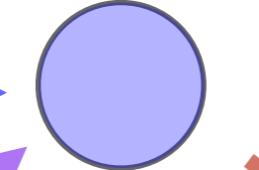
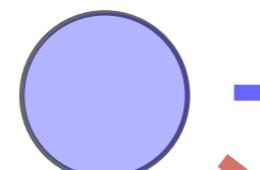
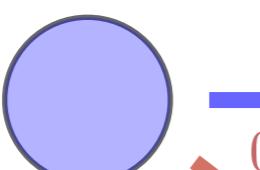
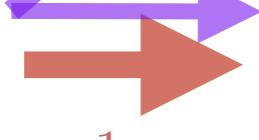
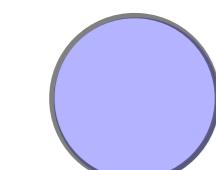
$t = 2$

$t = 3$

$t = 4$

$t = 5$

OK



$t = 0$

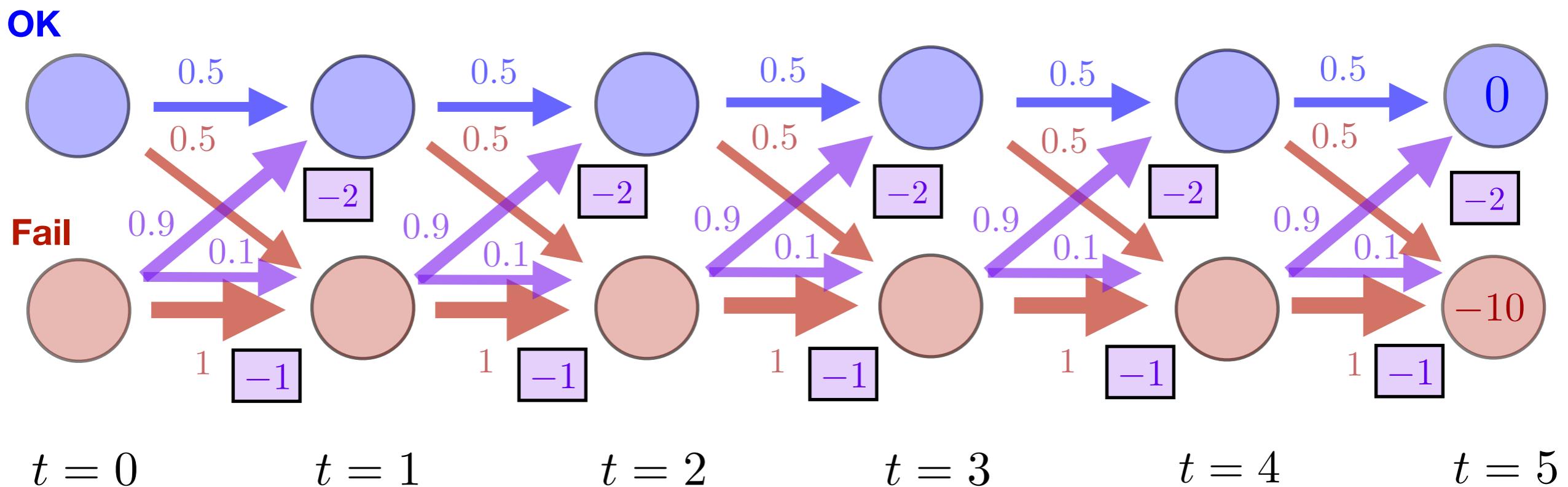
$t = 1$

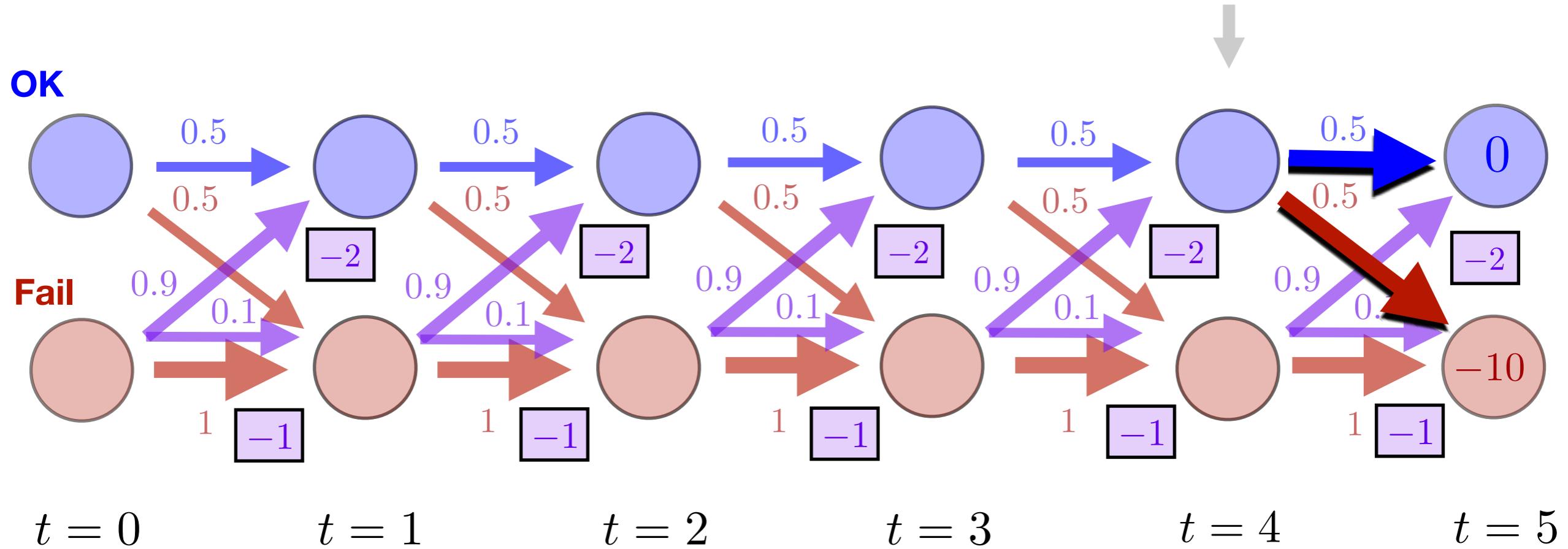
$t = 2$

$t = 3$

$t = 4$

$t = 5$

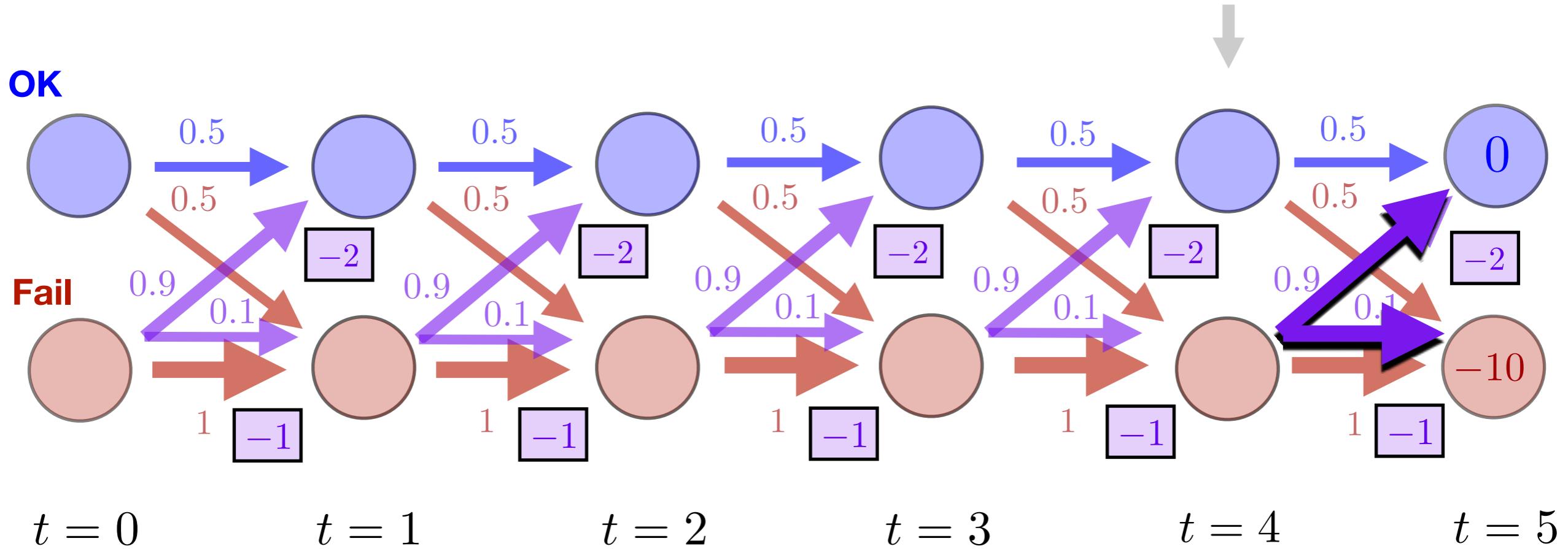




From State 1

action 1

$$-5 = 0 \times 0.5 + -10 \times 0.5$$



From State 1

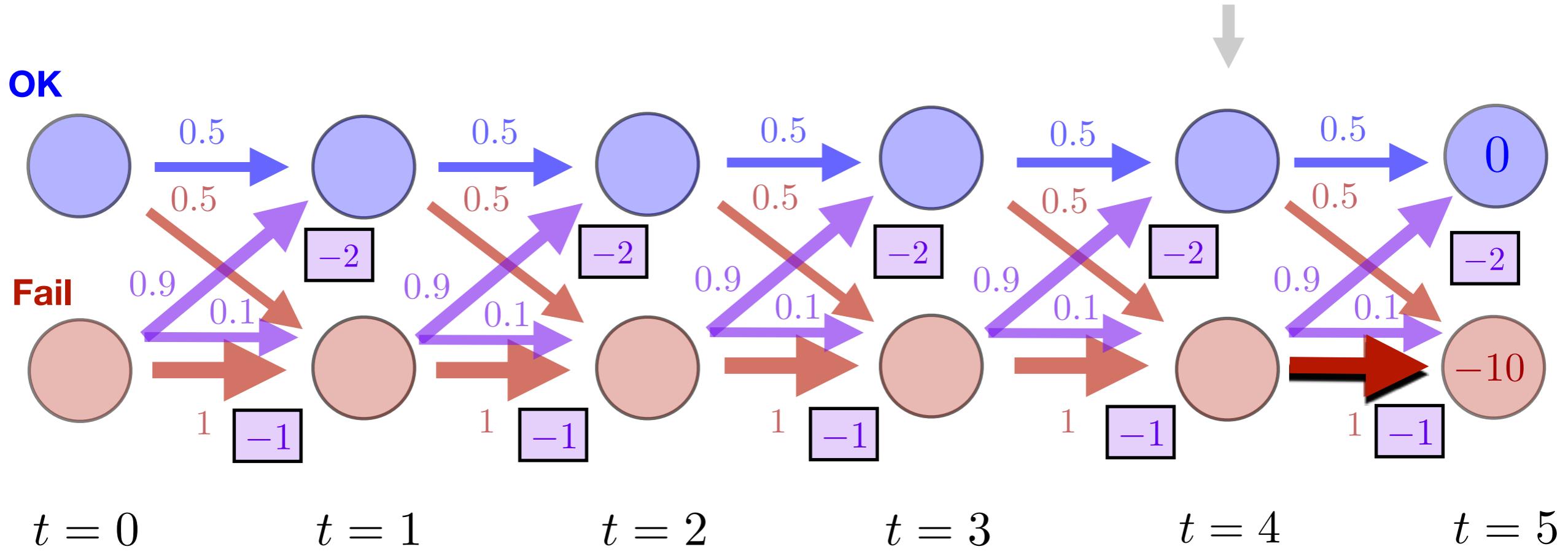
action 1

$$-5 = 0 \times 0.5 + -10 \times 0.5$$

From State 2

action 1

$$-3 = 0 \times 0.9 + -10 \times 0.1 + -2$$



From State 1

action 1

$$-5 = 0 \times 0.5 + -10 \times 0.5$$

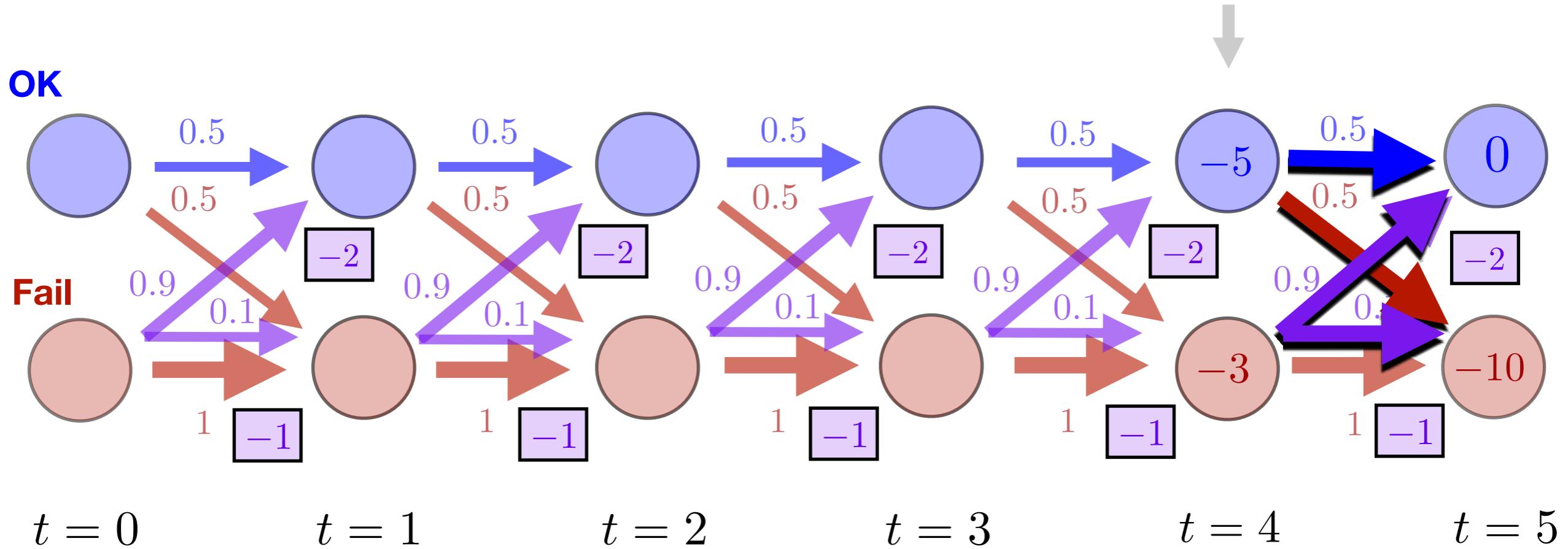
From State 2

action 1

$$-3 = 0 \times 0.9 + -10 \times 0.1 + -2$$

action 2

$$-11 = -10 \times 1 + -1$$



From State 1

action 1

$$-5 = \text{blue circle} \times 0.5 + \text{red circle} \times 0.5$$

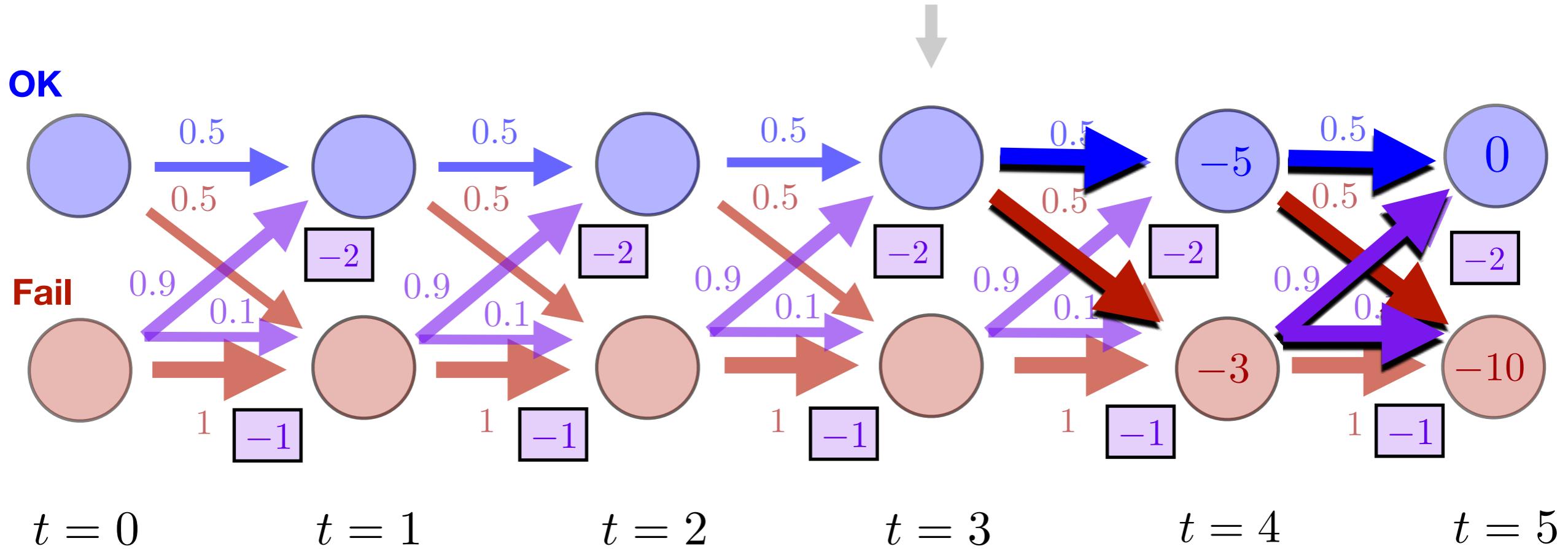
From State 2

action 1

$$-3 = \text{blue circle} \times 0.9 + \text{red circle} \times 0.1 + \text{purple box}$$

action 2

$$-11 = \text{red circle} \times 1 + \text{purple box}$$



From State 1

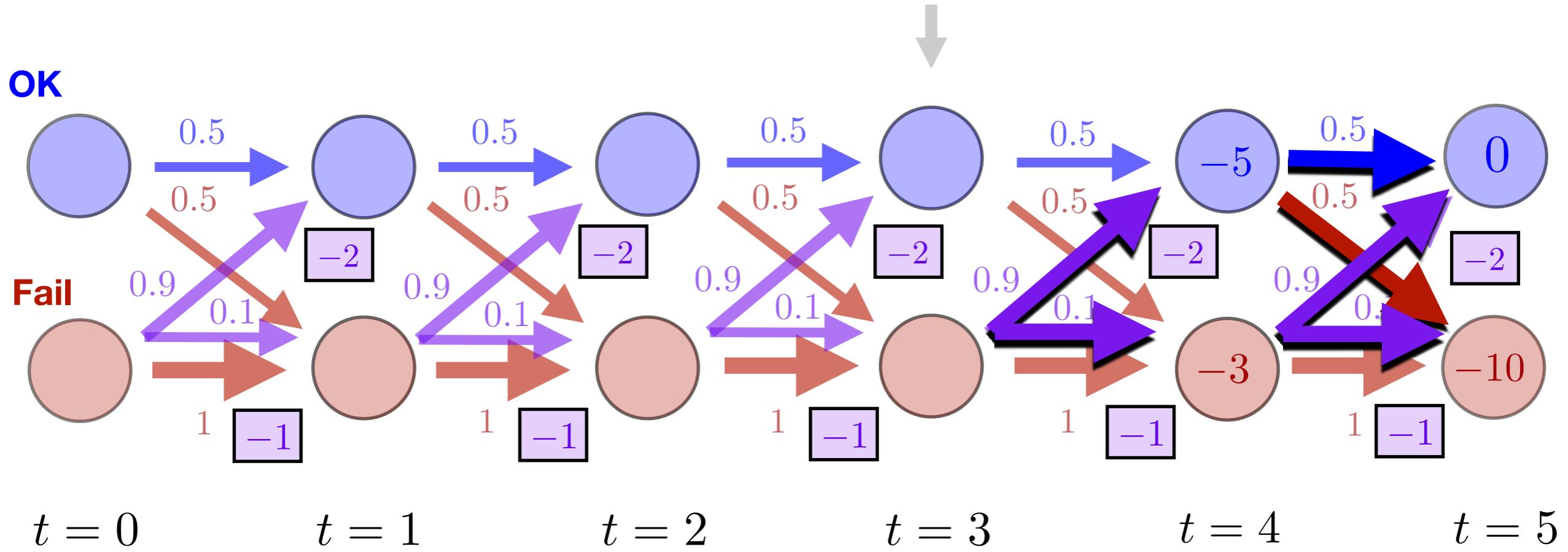
action 1

$$-4 = \text{blue circle} \times 0.5 + \text{red circle} \times 0.5$$

From State 2

action 1

action 2



From State 1

action 1

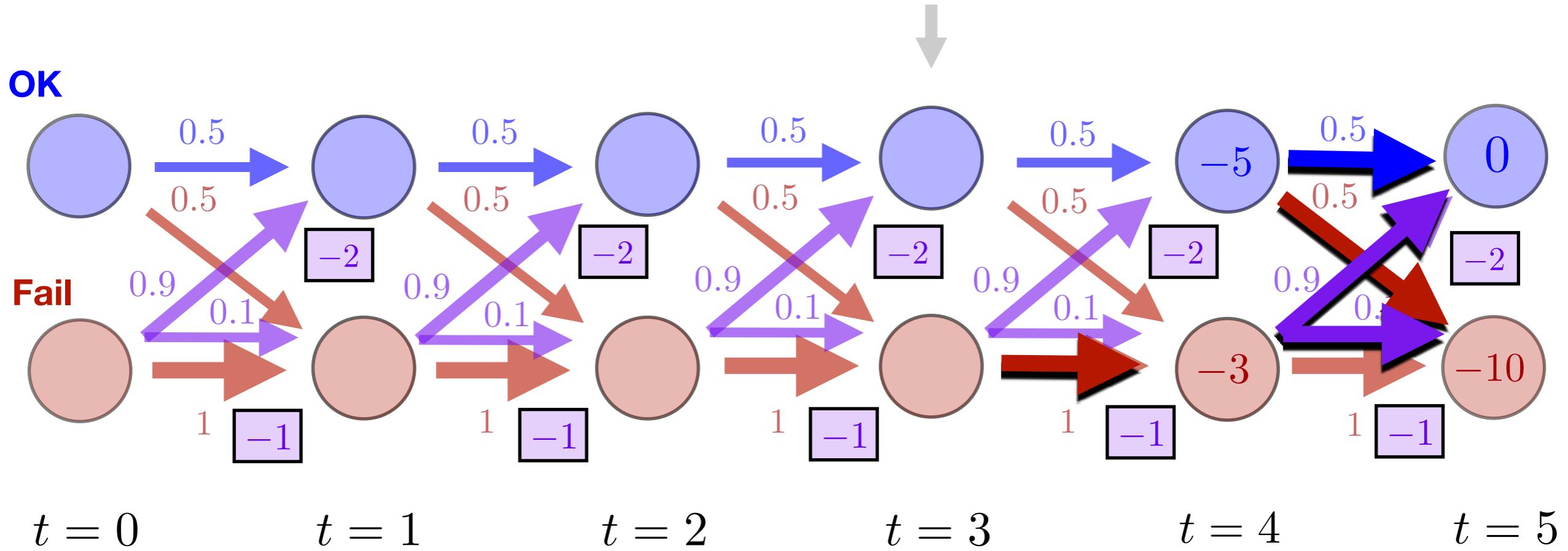
$$-4 = \text{blue circle} \times 0.5 + \text{red circle} \times 0.5$$

From State 2

action 1

$$-6.8 = \text{blue circle} \times 0.9 + \text{red circle} \times 0.1 + \text{purple square} \times 0.5$$

action 2



From State 1

action 1

$$-4 = \text{OK} \times 0.5 + \text{Fail} \times 0.5$$

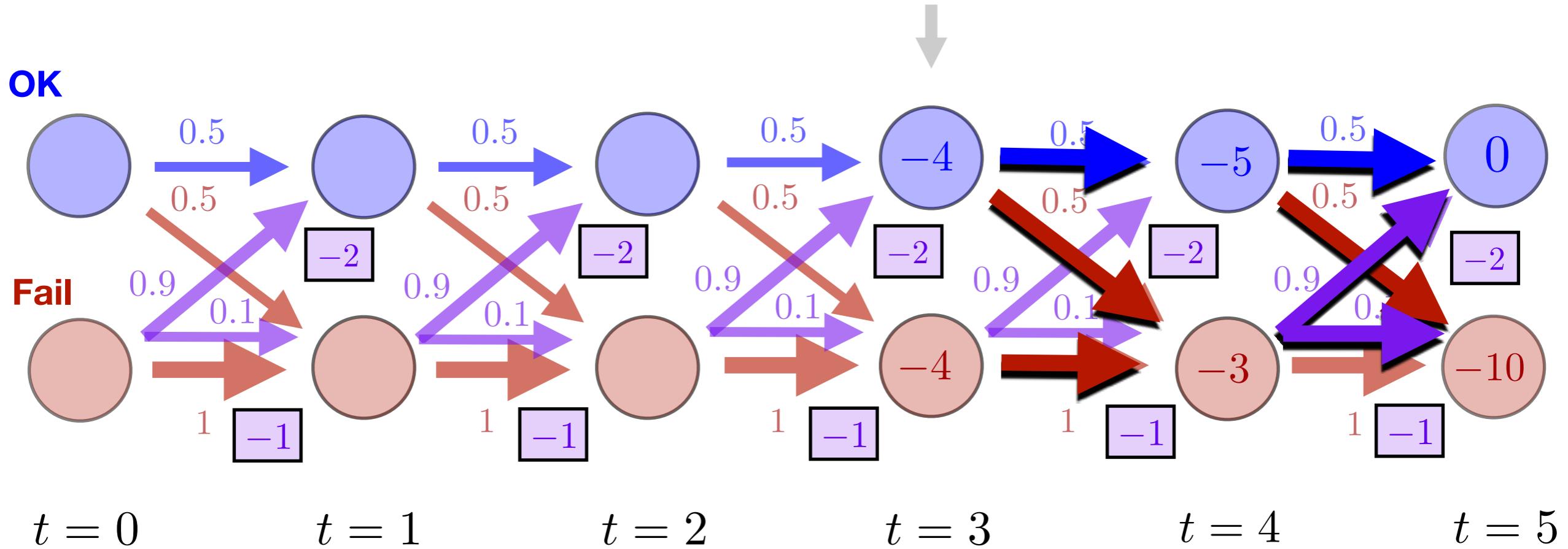
From State 2

action 1

$$-6.8 = \text{OK} \times 0.9 + \text{Fail} \times 0.1 + \text{Reward}$$

action 2

$$-4 = \text{Fail} \times 1 + \text{Reward}$$



From State 1

action 1

$$-4 = \text{blue circle} \times 0.5 + \text{red circle} \times 0.5$$

max

From State 2

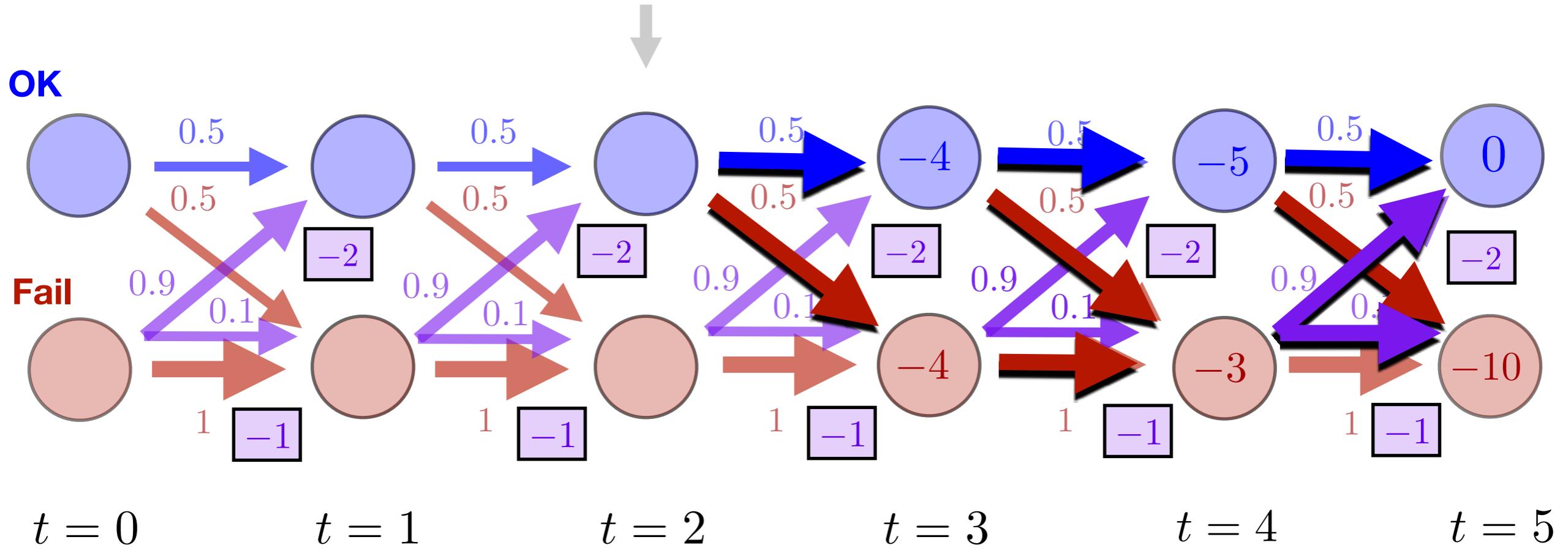
action 1

$$-6.8 = \text{blue circle} \times 0.9 + \text{red circle} \times 0.1 + \text{purple box}$$

max

action 2

$$-4 = \text{red circle} \times 1 + \text{purple box}$$



From State 1

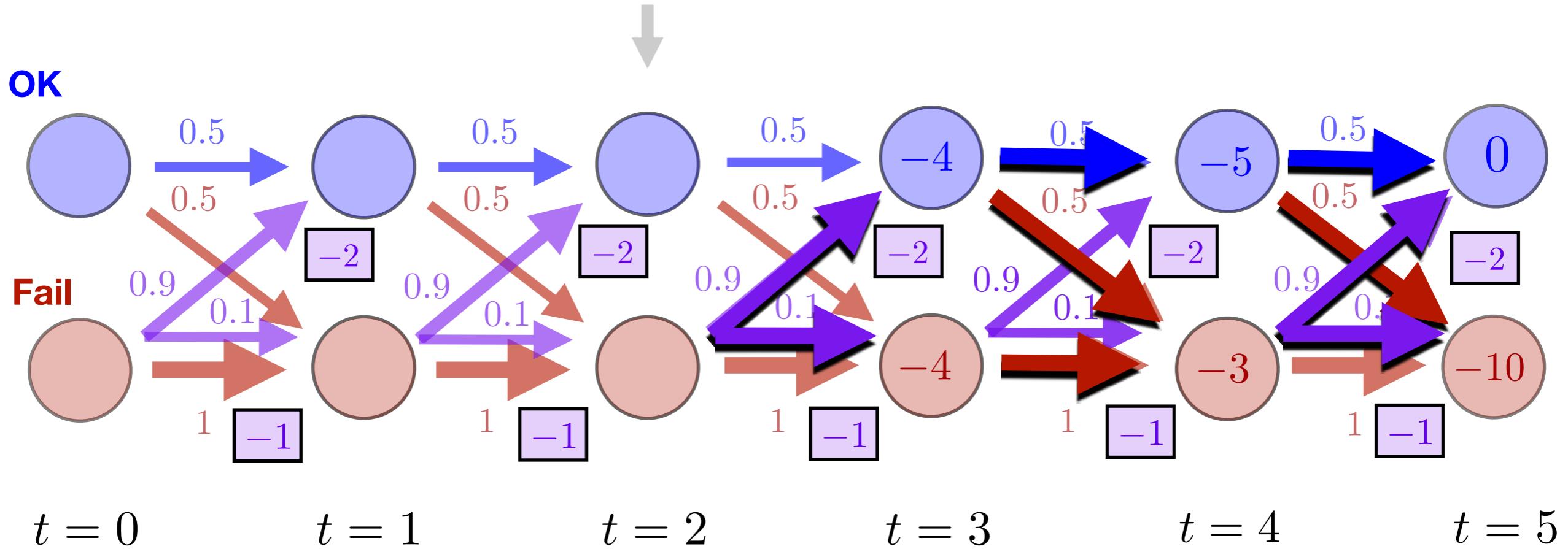
action 1

$$-4 = \text{blue circle} \times 0.5 + \text{red circle} \times 0.5$$

From State 2

action 1

action 2



From State 1

action 1

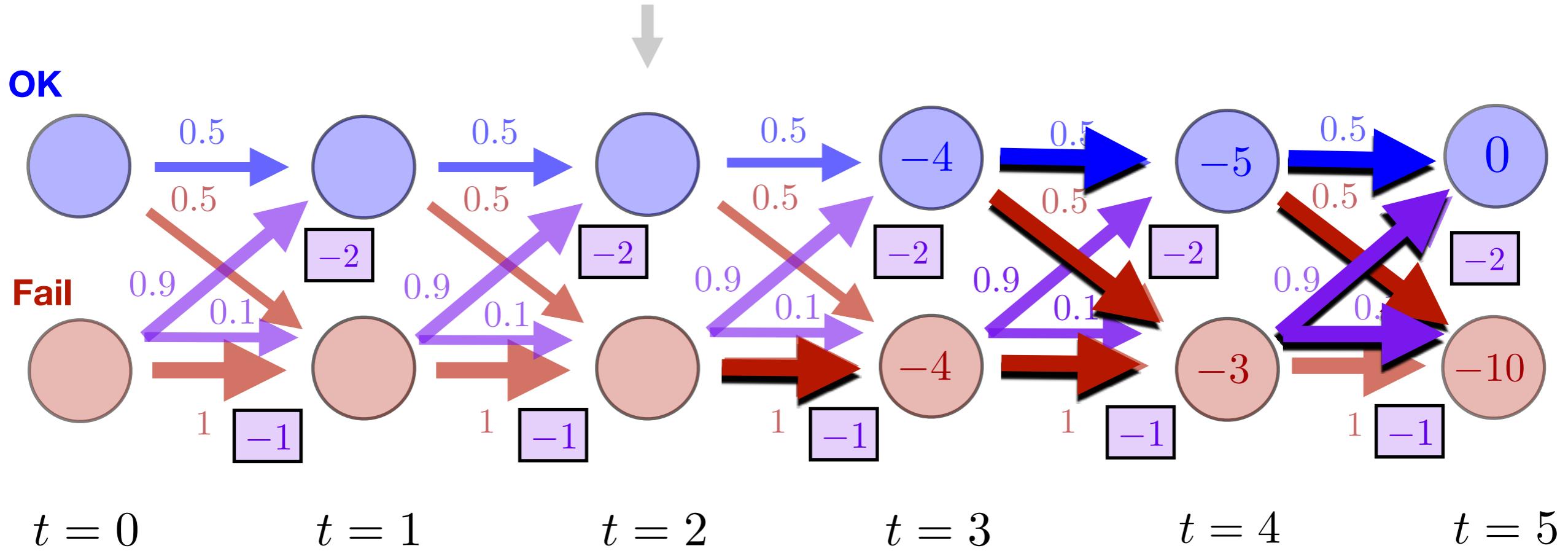
$$-4 = \text{blue circle} \times 0.5 + \text{red circle} \times 0.5$$

From State 2

action 1

$$-6 = \text{blue circle} \times 0.9 + \text{red circle} \times 0.1 + \text{purple box}$$

action 2



From State 1

action 1

$$-4 = \text{blue circle} \times 0.5 + \text{red circle} \times 0.5$$

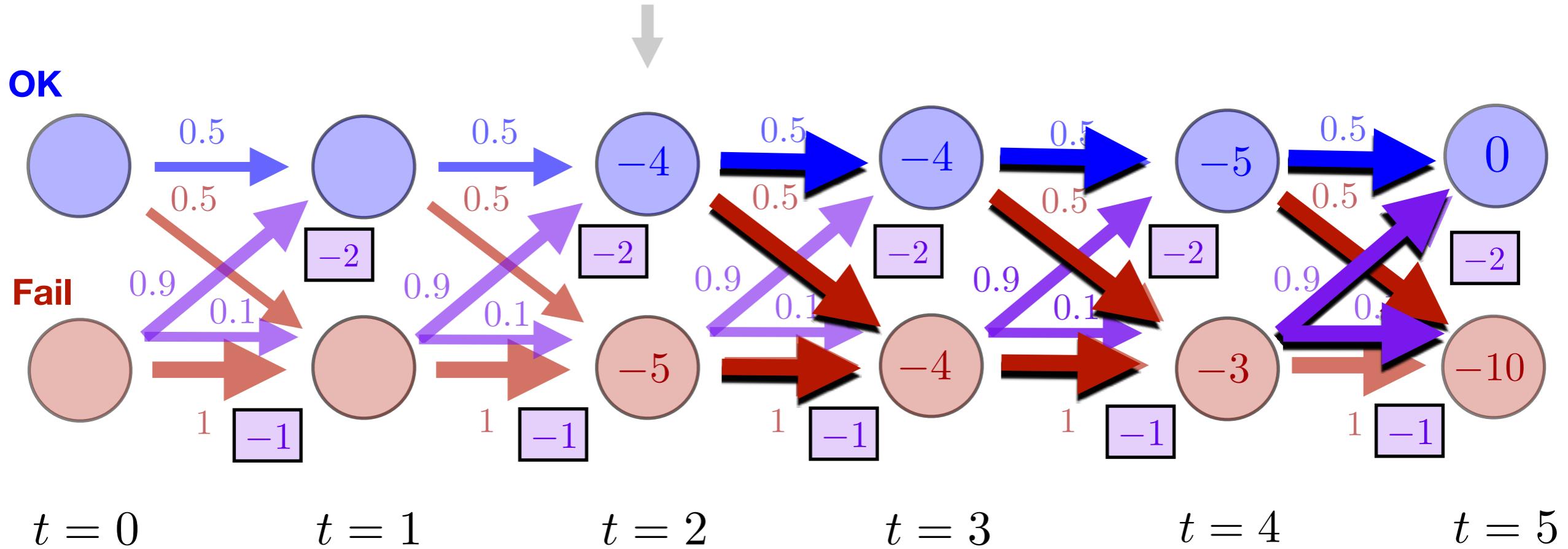
From State 2

action 1

$$-6 = \text{blue circle} \times 0.9 + \text{red circle} \times 0.1 + \text{purple box}$$

action 2

$$-5 = \text{red circle} \times 1 + \text{purple box}$$



From State 1

action 1

$$-4 = \text{blue circle} \times 0.5 + \text{red circle} \times 0.5$$

max

From State 2

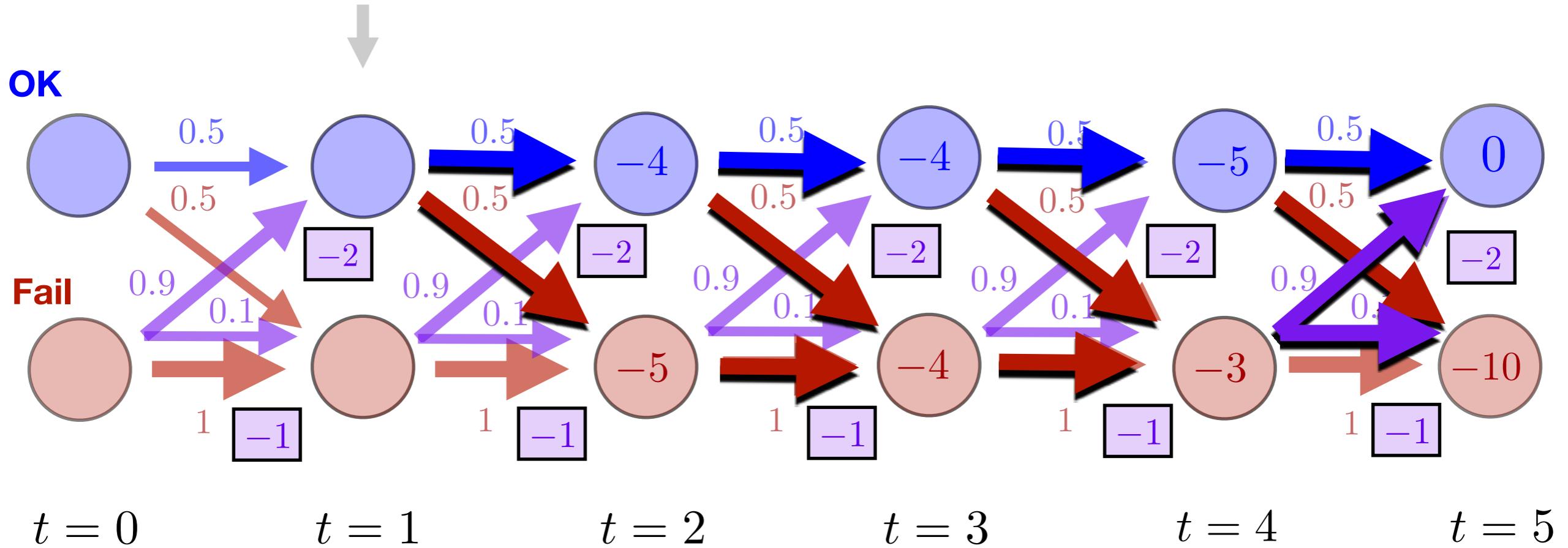
action 1

$$-6 = \text{blue circle} \times 0.9 + \text{red circle} \times 0.1 + \text{purple box} \times 0.0$$

max

action 2

$$-5 = \text{red circle} \times 1 + \text{purple box} \times 0.0$$



From State 1

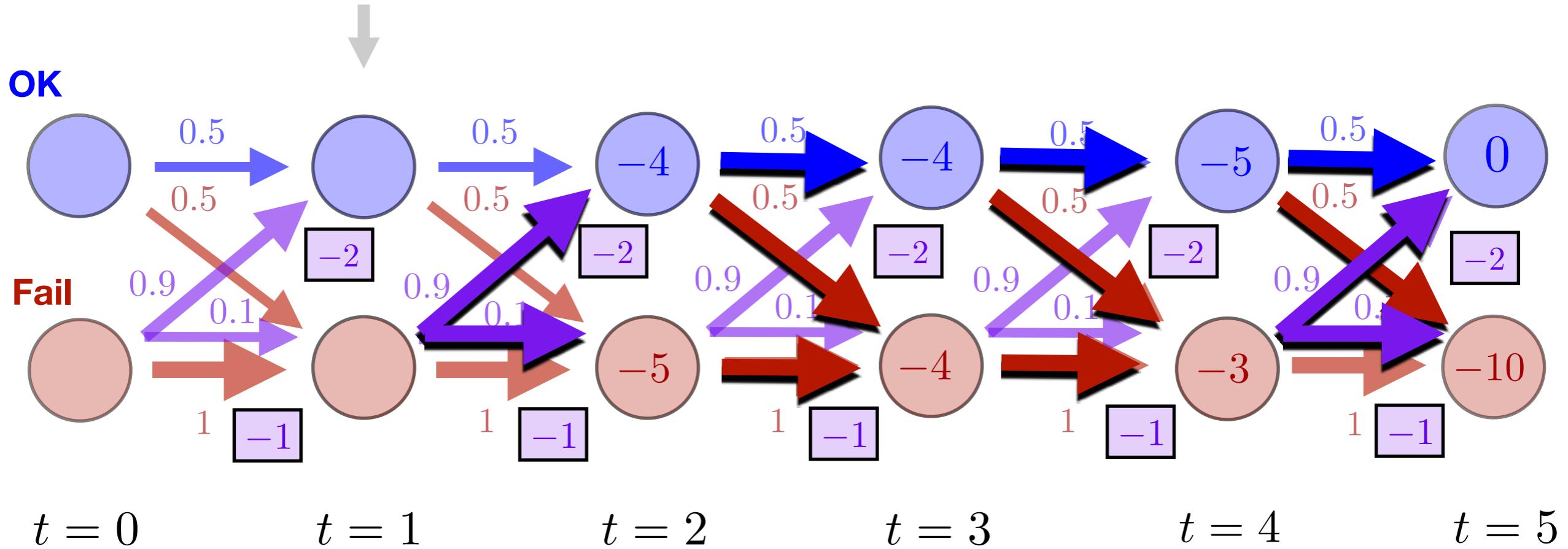
action 1

$$-4.5 = \text{blue circle} \times 0.5 + \text{red circle} \times 0.5$$

From State 2

action 1

action 2



From State 1

action 1

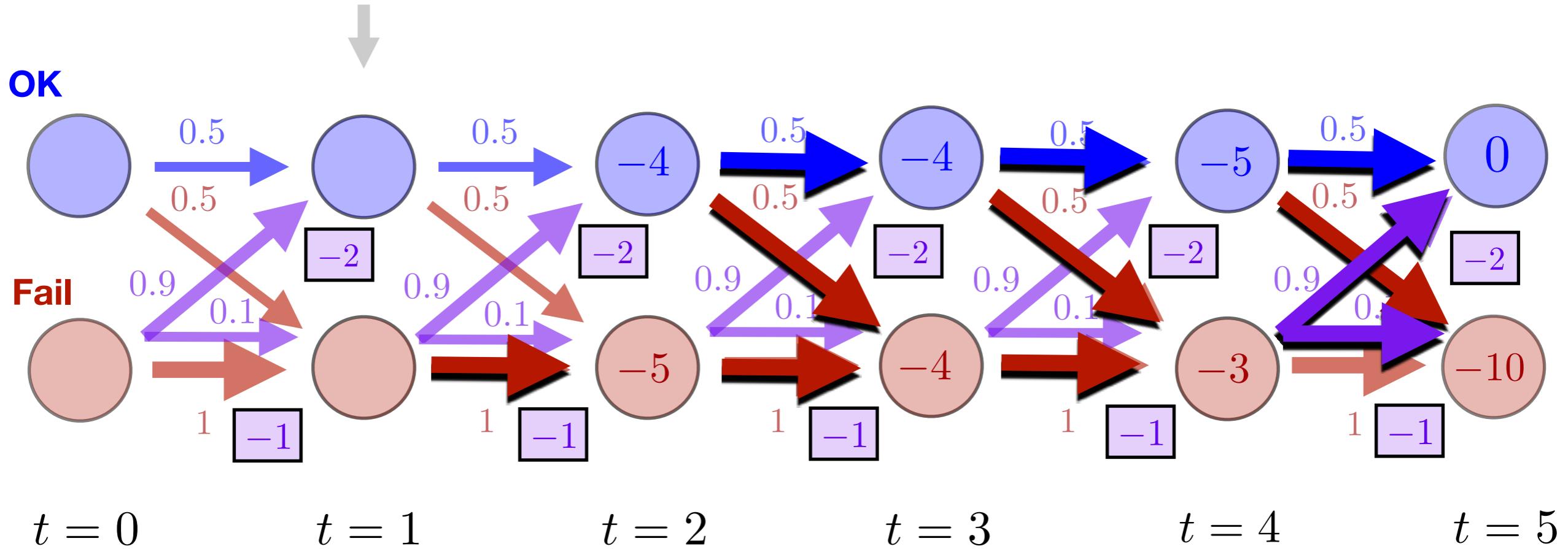
$$-4.5 = \text{---4} \times 0.5 + \text{---5} \times 0.5$$

From State 2

action 1

$$-6.1 = \text{---4} \times 0.9 + \text{---5} \times 0.1 + \text{---2}$$

action 2



From State 1

action 1

$$-4.5 = \text{blue circle} \times 0.5 + \text{red circle} \times 0.5$$

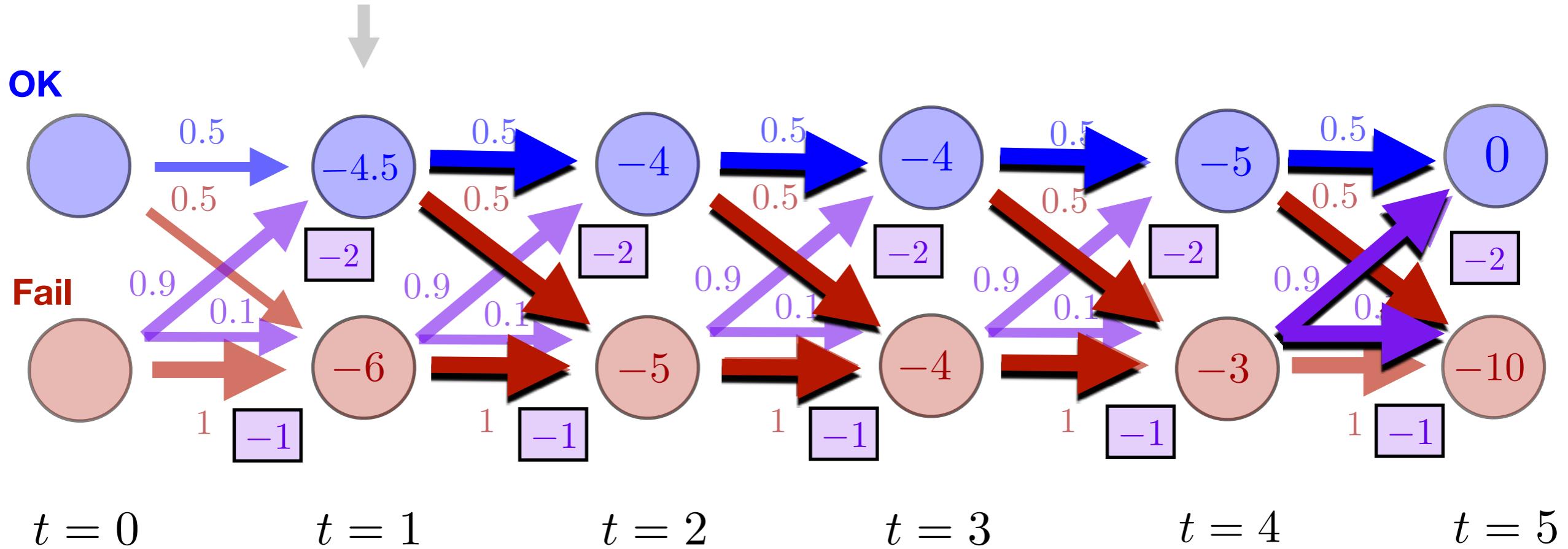
From State 2

action 1

$$-6.1 = \text{blue circle} \times 0.9 + \text{red circle} \times 0.1 + \text{purple square}$$

action 2

$$-6 = \text{red circle} \times 1 + \text{purple square}$$



From State 1

action 1

$$-4.5 = \underset{\max}{\text{---}} \times 0.5 + \underset{\max}{\text{---}} \times 0.5$$

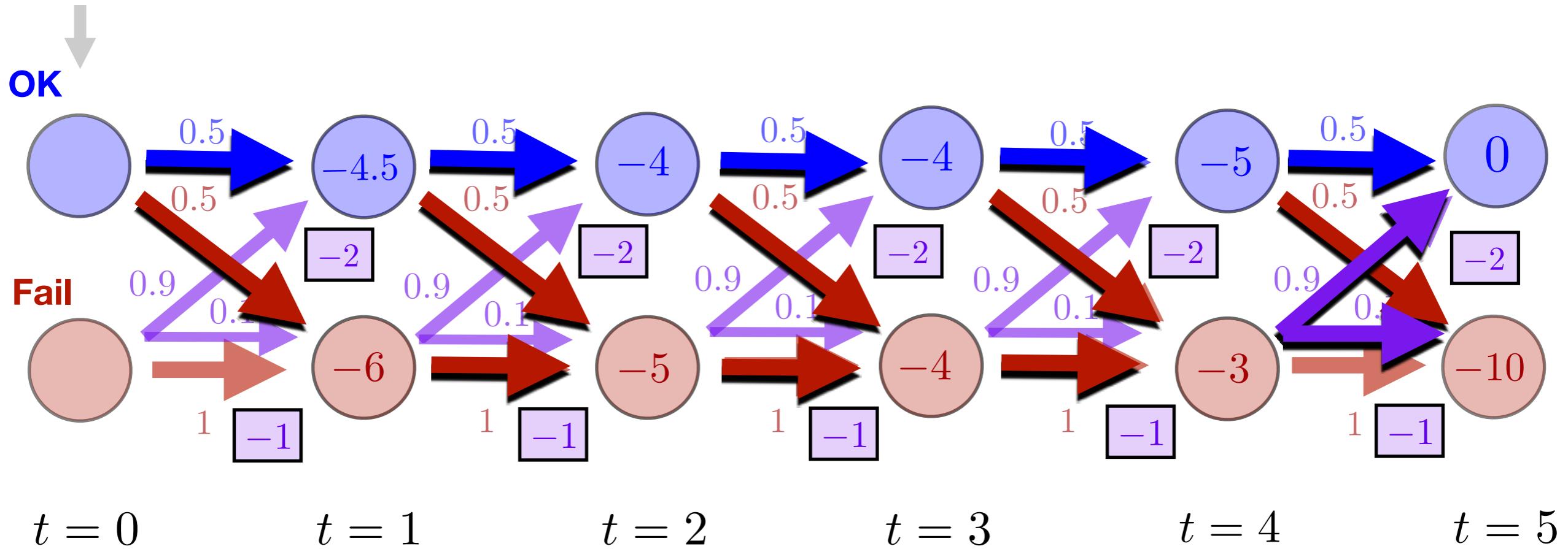
From State 2

action 1

$$-6.1 = \underset{\max}{\text{---}} \times 0.9 + \underset{\max}{\text{---}} \times 0.1 + \underset{\max}{\text{---}}$$

action 2

$$-6 = \underset{\max}{\text{---}} \times 1 + \underset{\max}{\text{---}}$$



From State 1

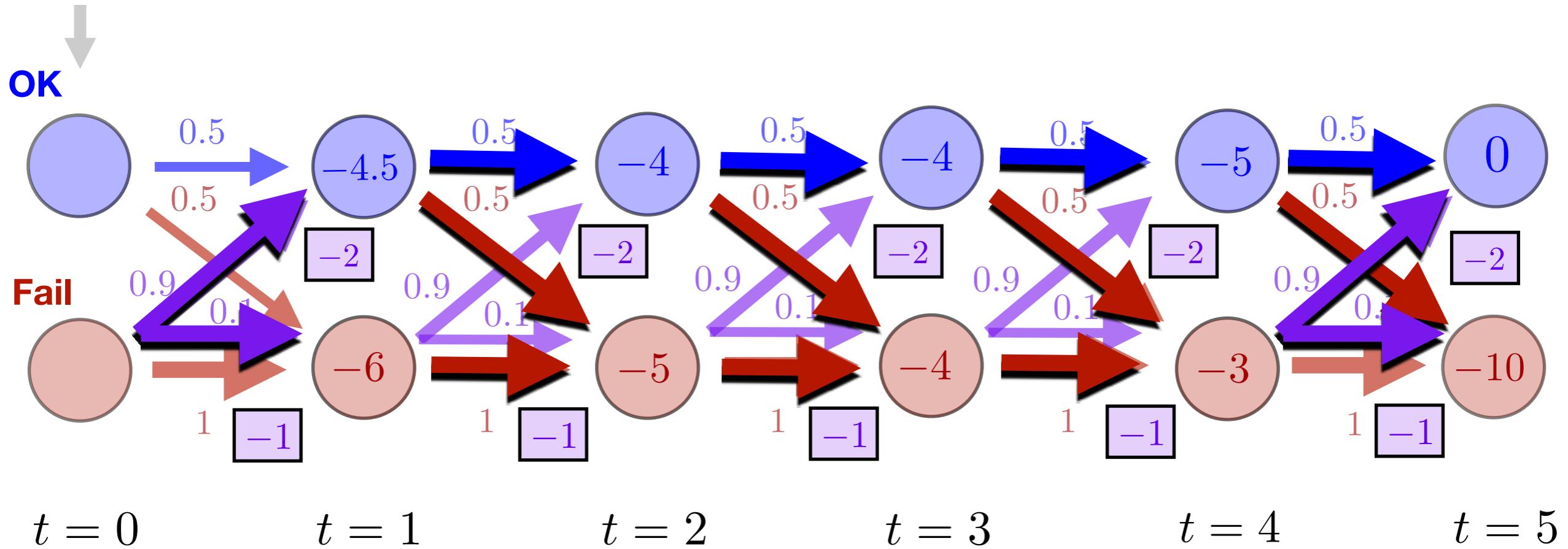
action 1

$$-5.3 = -4.5 \times 0.5 + -6 \times 0.5$$

From State 2

action 1

action 2



From State 1

action 1

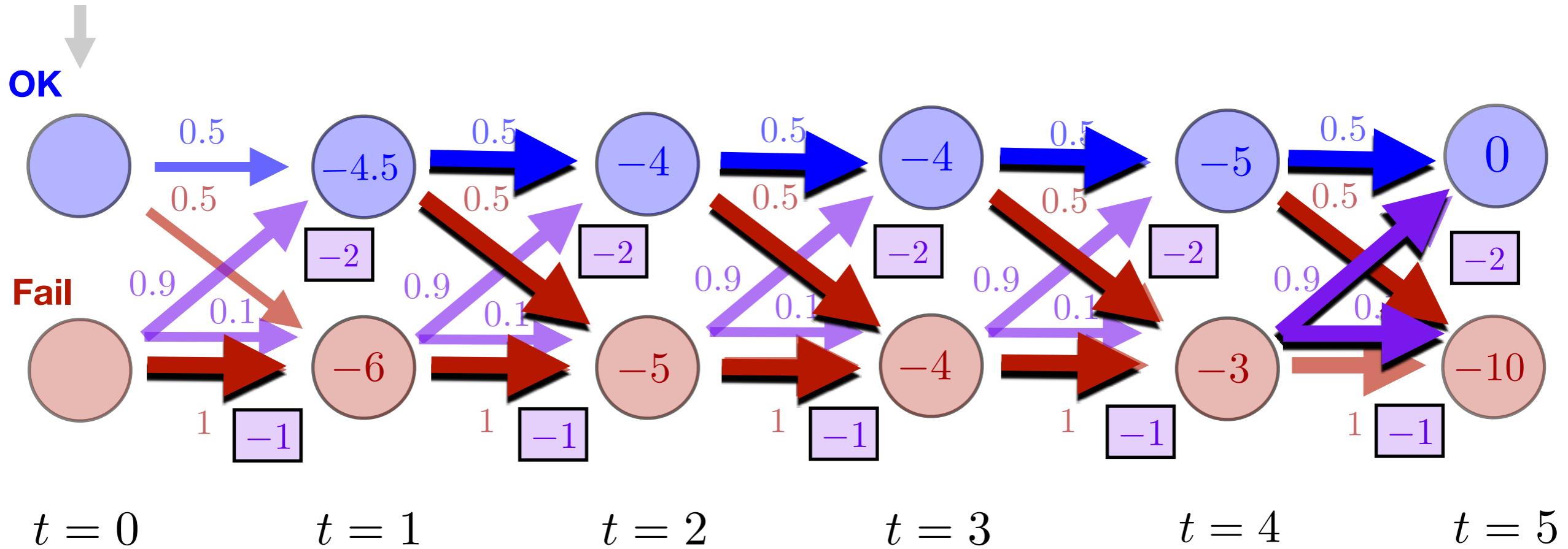
$$-5.3 = -4.5 \times 0.5 + -6 \times 0.5$$

From State 2

action 1

$$-6.7 = -4.5 \times 0.9 + -6 \times 0.1 + -2$$

action 2



From State 1

action 1

$$-5.3 = -4.5 \times 0.5 + -6 \times 0.5$$

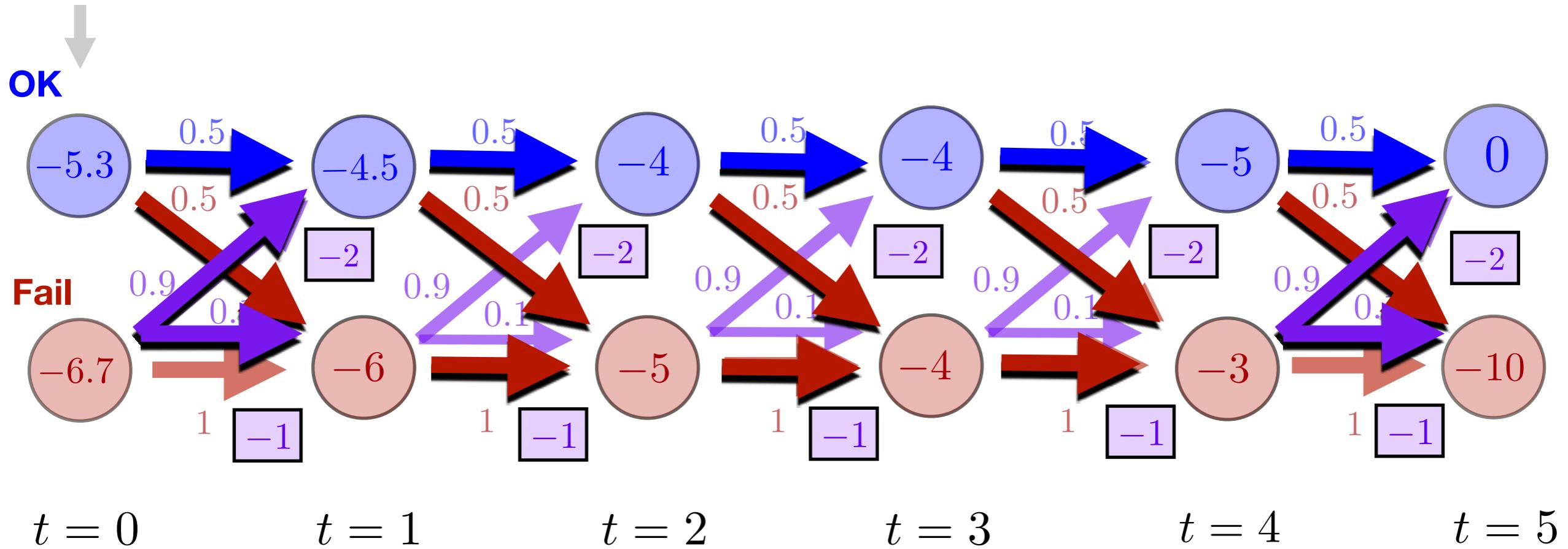
From State 2

action 1

$$-6.7 = -4.5 \times 0.9 + -6 \times 0.1 + -2$$

action 2

$$-7 = -6 \times 1 + -1$$



From State 1

action 1

$$-5.3 = \max(-4.5 \times 0.5 + -6 \times 0.5)$$

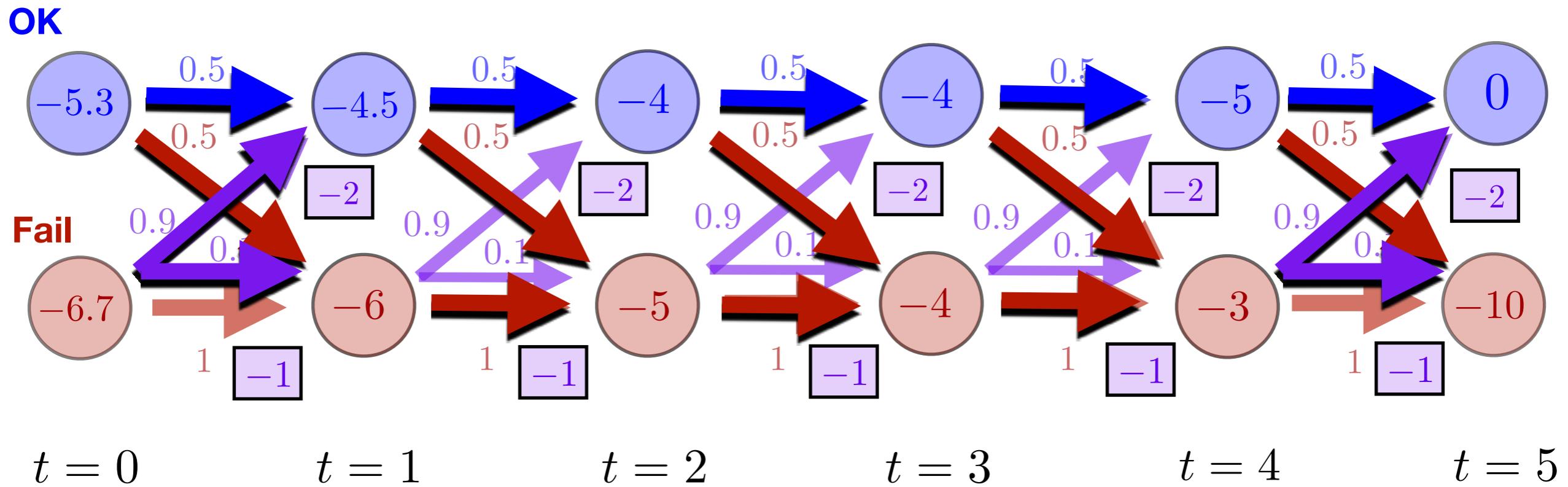
From State 2

action 1

$$-6.7 = \max(-4.5 \times 0.9 + -6 \times 0.1 + -2)$$

action 2

$$-7 = \max(-6 \times 1 + -1)$$



Action-Value

from state s , action a

$$q_{sa}(k-1) = \sum_{s'} v_{s'}(k) P(s'|s, a) + r_{sa}(k)$$

*"q-value"
general
form*

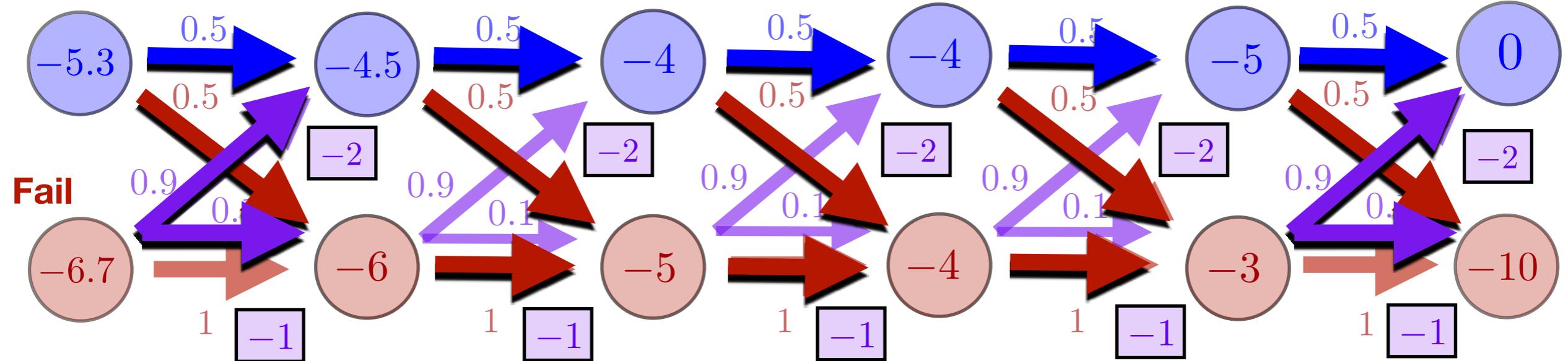
From State 2

action 1

$$-6.7 = -4.5 \times 0.9 + -6 \times 0.1 + -2$$

example

OK



$t = 0$

$t = 1$

$t = 2$

$t = 3$

$t = 4$

$t = 5$

Action-Value

from state s , action a

$$q_{sa}(k-1) = \sum_{s'} v_{s'}(k) P(s'|s, a) + r_{sa}(k)$$

*"q-value"
general
form*

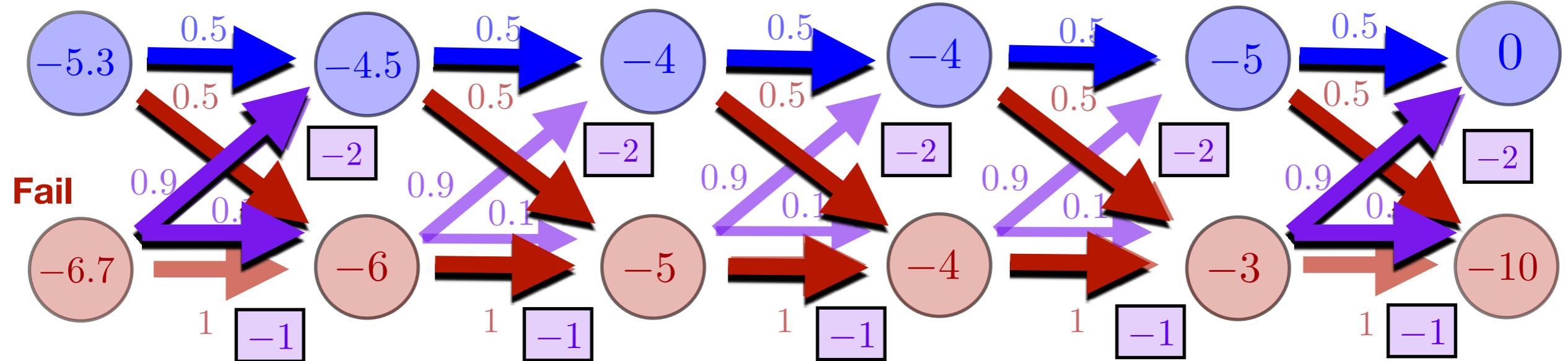
State-Value

for each state s

$$v_s(k-1) = \max_{a \in \mathcal{A}_s} q_{sa}(k-1)$$

"Reward-to-go"

OK



$t = 0$

$t = 1$

$t = 2$

$t = 3$

$t = 4$

$t = 5$

Action-Value

from state s , action a

$$q_{sa}(k-1) = \sum_{s'} v_{s'}(k) P(s'|s, a) + r_{sa}(k)$$

*"q-value"
general
form*

State-Value

for each state s

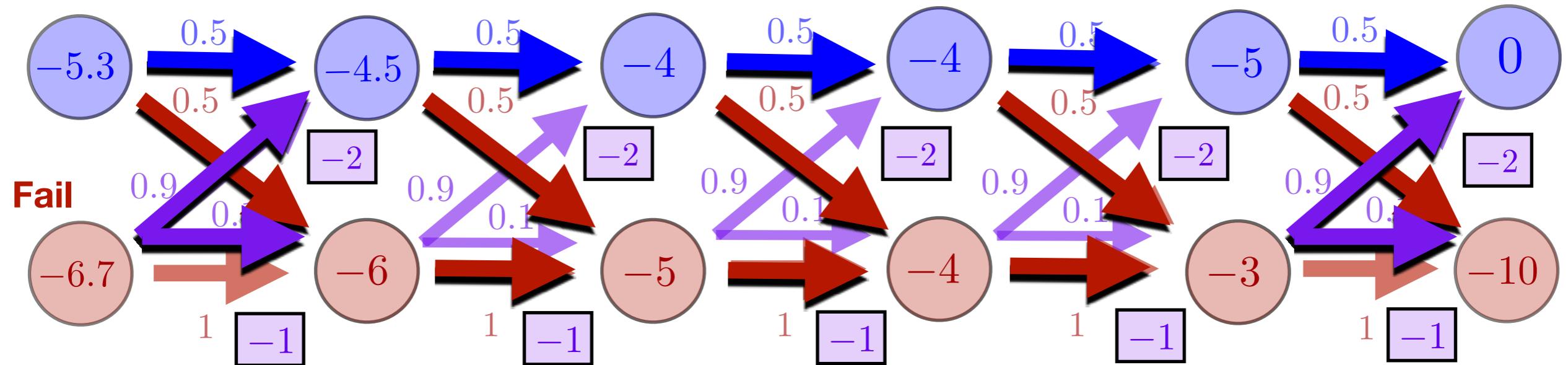
$$v_s(k-1) = \max_{a \in \mathcal{A}_s} q_{sa}(k-1)$$

"Reward-to-go"

Bellman Equation
Dynamic Programming

$$= \max_{a \in \mathcal{A}_s} \sum_{s'} v_{s'}(k) P(s'|s, a) + r_{sa}(k)$$

OK



$t = 0$

$t = 1$

$t = 2$

$t = 3$

$t = 4$

$t = 5$

State-Value Meaning

from state s

$v_s(k)$ = expected reward till the end of the time horizon if you start in state s at time k and use optimal policy

State-Value

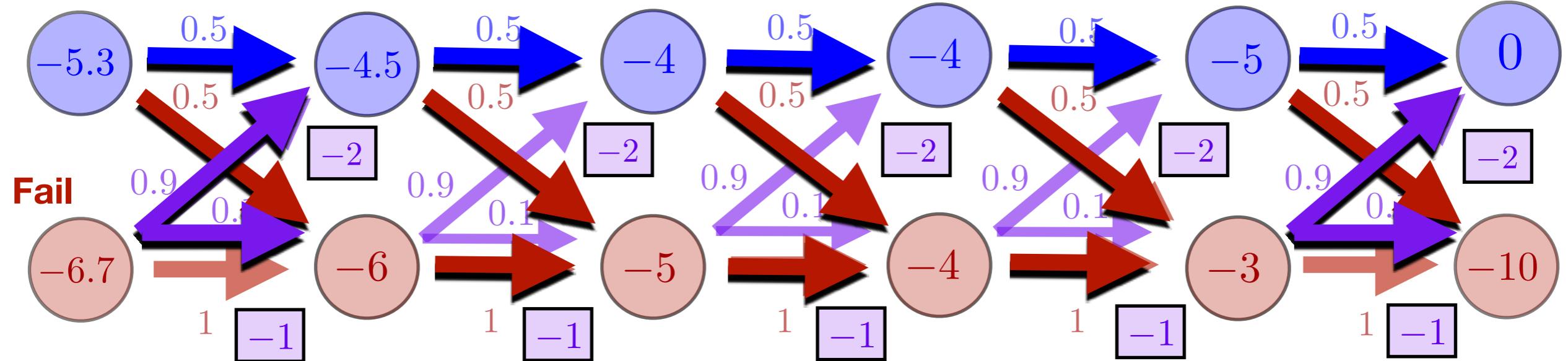
for each state s

$$v_s(k-1) = \max_{a \in \mathcal{A}_s} q_{sa}(k-1) \quad \text{"Reward-to-go"}$$

Bellman Equation
Dynamic Programming

$$= \max_{a \in \mathcal{A}_s} \sum_{s'} v_{s'}(k) P(s'|s, a) + r_{sa}(k)$$

OK



$t = 0$

$t = 1$

$t = 2$

$t = 3$

$t = 4$

$t = 5$

State-Value Meaning

from state s

$v_s(0)$ = expected reward if you start in state s (at time 0)
and use the optimal policy

State-Value

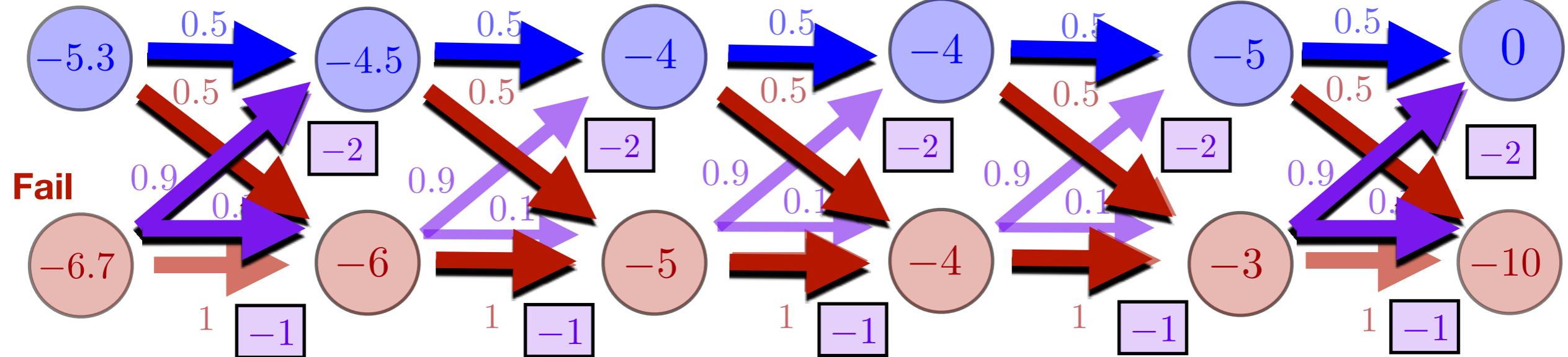
for each state s

$v_s(k-1) = \max_{a \in \mathcal{A}_s} q_{sa}(k-1)$ “Reward-to-go”

Bellman Equation
Dynamic Programming

$$= \max_{a \in \mathcal{A}_s} \sum_{s'} v_{s'}(k) P(s'|s, a) + r_{sa}(k)$$

OK



$t = 0$

$t = 1$

$t = 2$

$t = 3$

$t = 4$

$t = 5$

Total Reward

(Finite Horizon)

$$R = \sum_s p_s(0) v_s(0)$$

State-Value

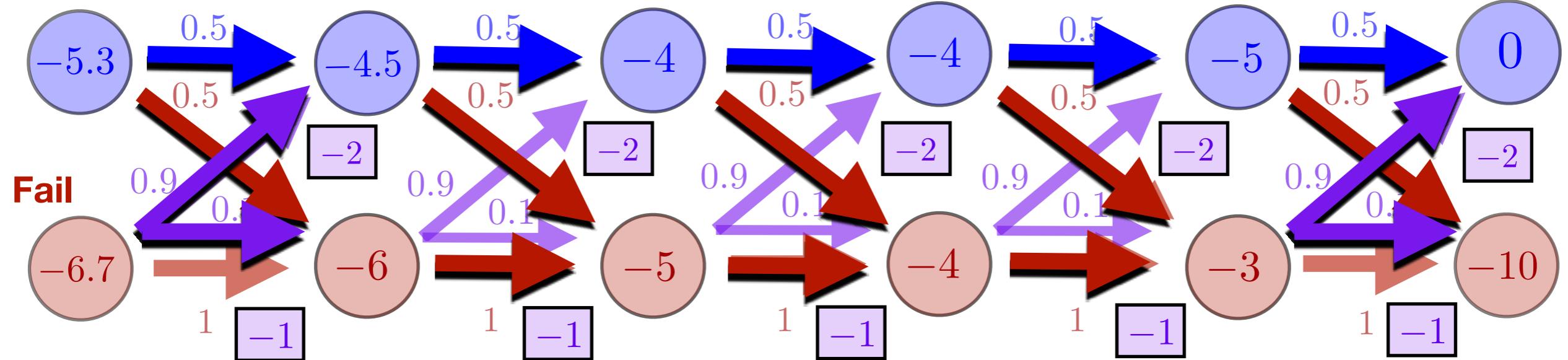
for each state s

$$v_s(k-1) = \max_{a \in \mathcal{A}_s} q_{sa}(k-1) \quad \text{"Reward-to-go"}$$

Bellman Equation
Dynamic Programming

$$= \max_{a \in \mathcal{A}_s} \sum_{s'} v_{s'}(k) P(s'|s, a) + r_{sa}(k)$$

OK



$t = 0$

$t = 1$

$t = 2$

$t = 3$

$t = 4$

$t = 5$

Total Reward
(Finite Horizon)

$$R = \sum_s p_s(0) \sum_{k=0}^{K-1} \sum_{s'} r_{s'a} p_{s'a}(k)$$

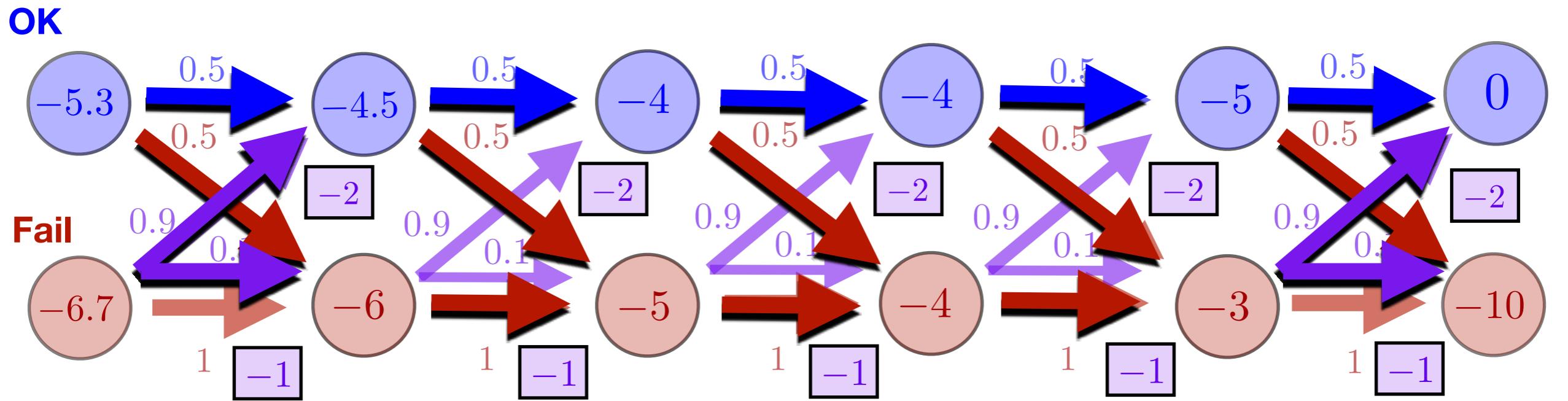
State-Value

for each state s

$$v_s(k-1) = \max_{a \in \mathcal{A}_s} q_{sa}(k-1) \quad \text{"Reward-to-go"}$$

Bellman Equation
Dynamic Programming

$$= \max_{a \in \mathcal{A}_s} \sum_{s'} v_{s'}(k) P(s'|s, a) + r_{sa}(k)$$



$t = 0$

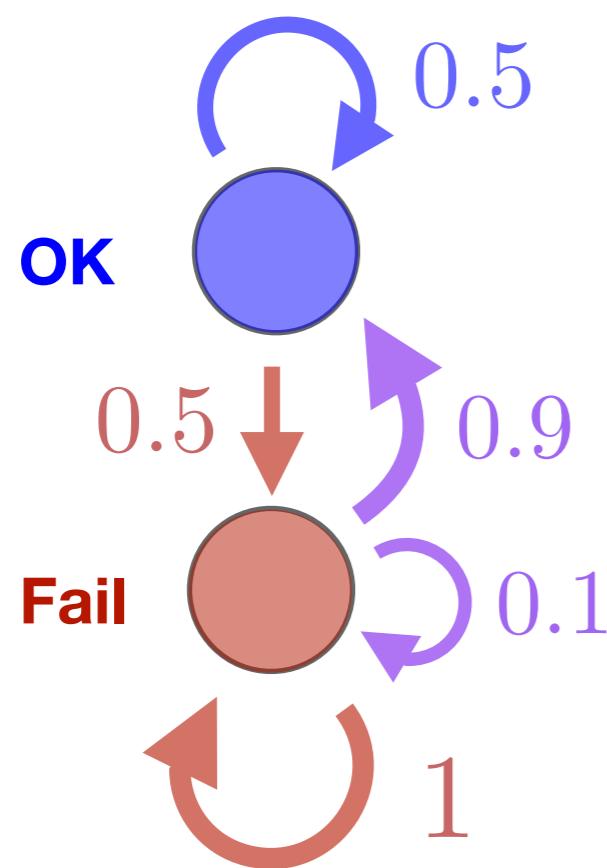
$t = 1$

$t = 2$

$t = 3$

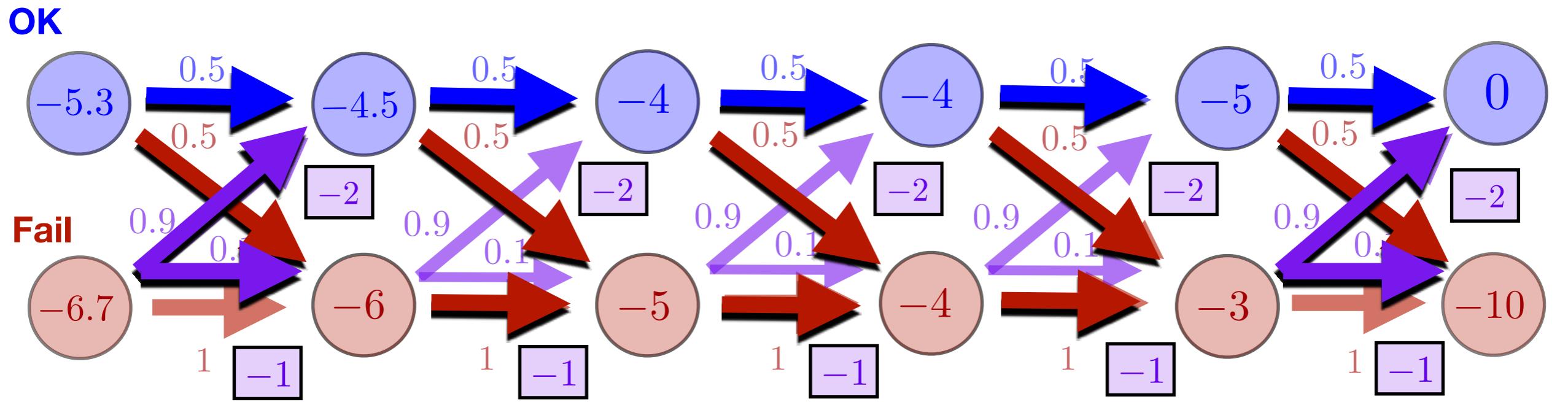
$t = 4$

$t = 5$



Bellman Equation

$$v_s(k-1) = \max_{a \in \mathcal{A}_s} \sum_{s'} v_{s'}(k) P(s'|s, a) + r_{sa}(k)$$



$t = 0$

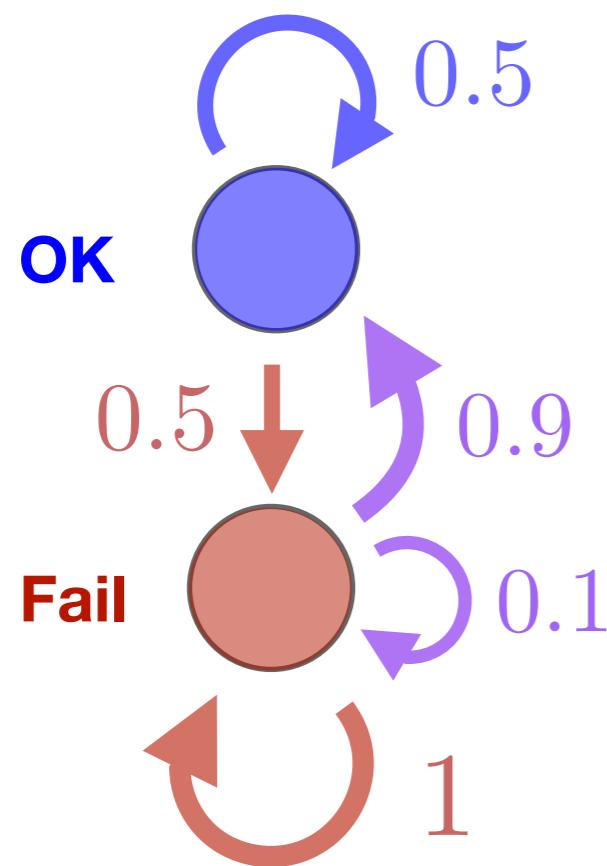
$t = 1$

$t = 2$

$t = 3$

$t = 4$

$t = 5$



Bellman Equation

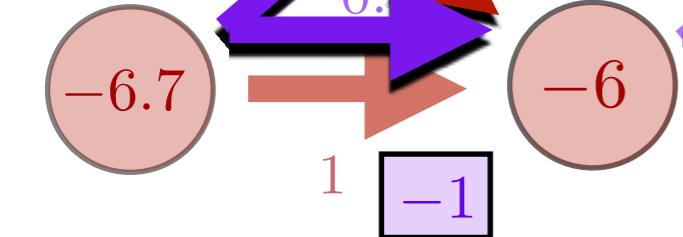
$$v_s(k-1) = \max_{a \in \mathcal{A}_s} \sum_{s'} v_{s'}(k) P(s'|s, a) + r_{sa}(k)$$

*iterative equation...
value doesn't necessarily converge...*

OK



Fail



$t = 0$

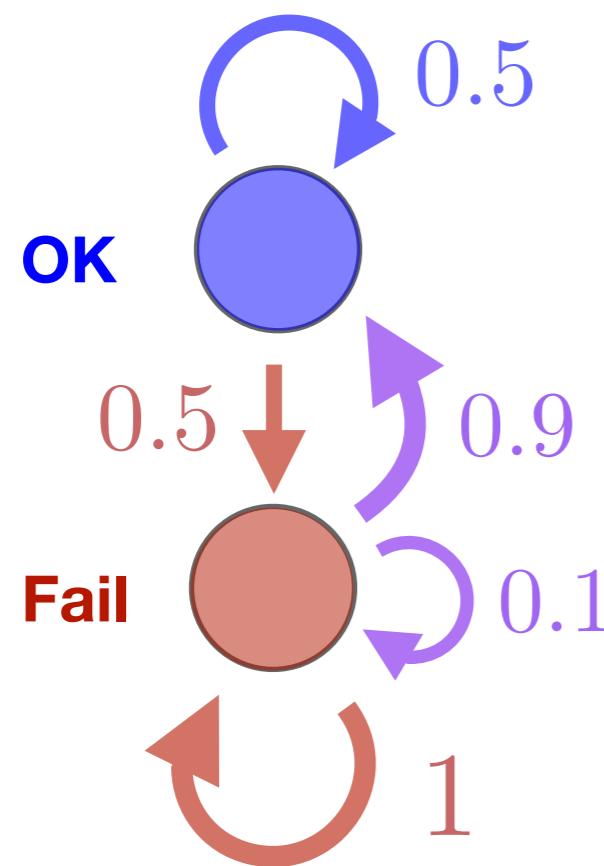
$t = 1$

$t = 2$

$t = 3$

$t = 4$

$t = 5$



Discounted Bellman Equation

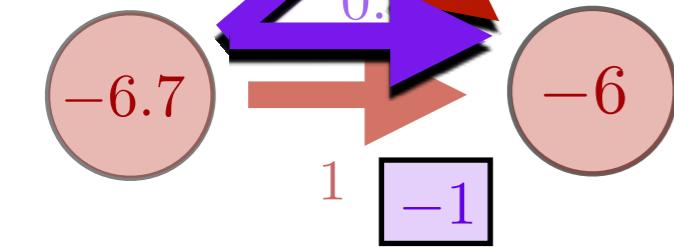
$$v_s(k-1) = \max_{a \in \mathcal{A}_s} \sum_{s'} \gamma v_{s'}(k) P(s'|s, a) + r_{sa}(k)$$

$\gamma \in [0, 1)$

OK



Fail



$t = 0$

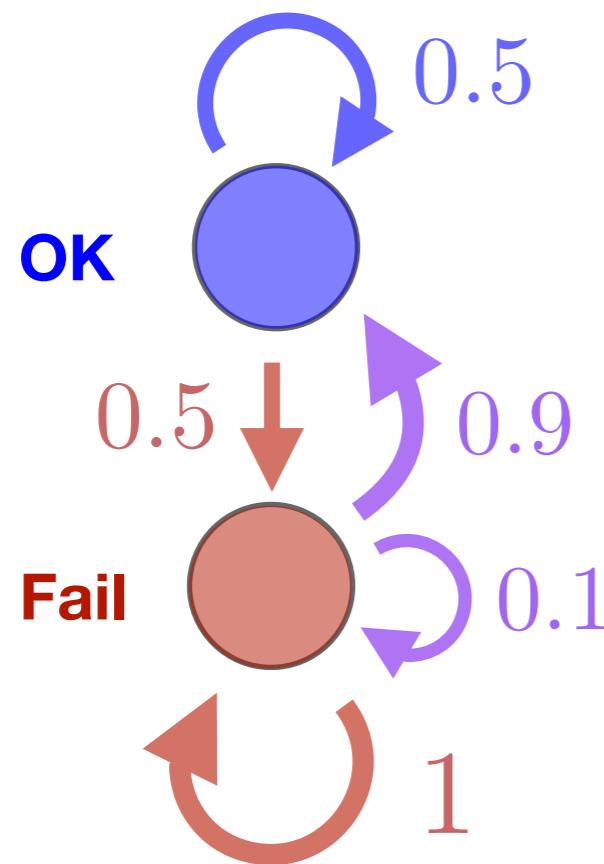
$t = 1$

$t = 2$

$t = 3$

$t = 4$

$t = 5$



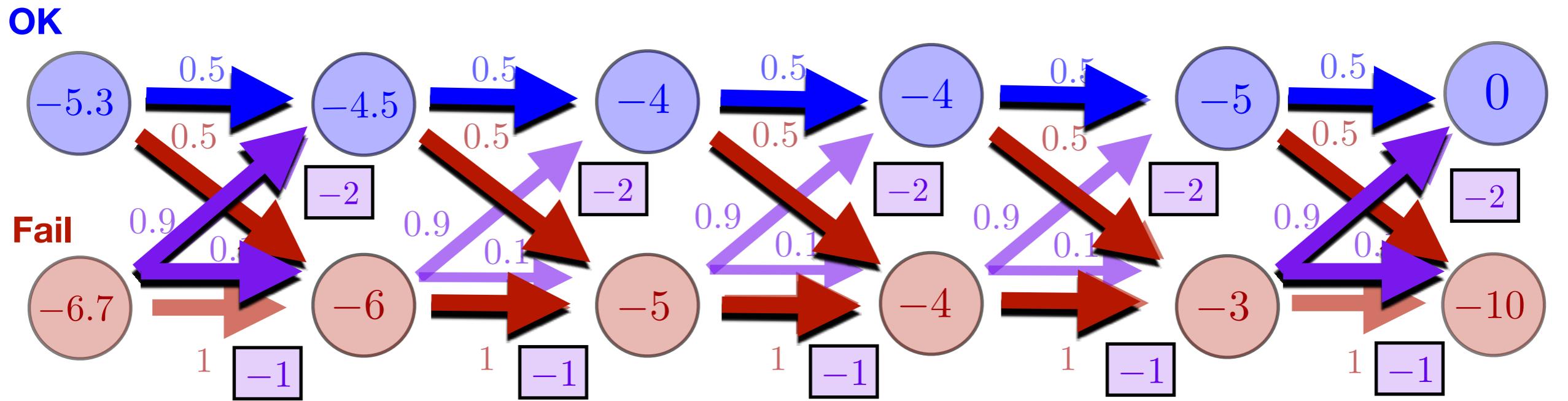
Discounted Bellman Equation

$$\gamma \in [0, 1)$$

$$v_s(k-1) = \max_{a \in \mathcal{A}_s} \sum_{s'} \gamma v_{s'}(k) P(s'|s, a) + r_{sa}(k)$$

if $0 \leq \gamma < 1$ we can always find $v_s(\infty)$ such that

$$v_s(\infty) = \max_{a \in \mathcal{A}_s} \sum_{s'} \gamma v_{s'}(\infty) P(s'|s, a) + r_{sa}$$



$t = 0$

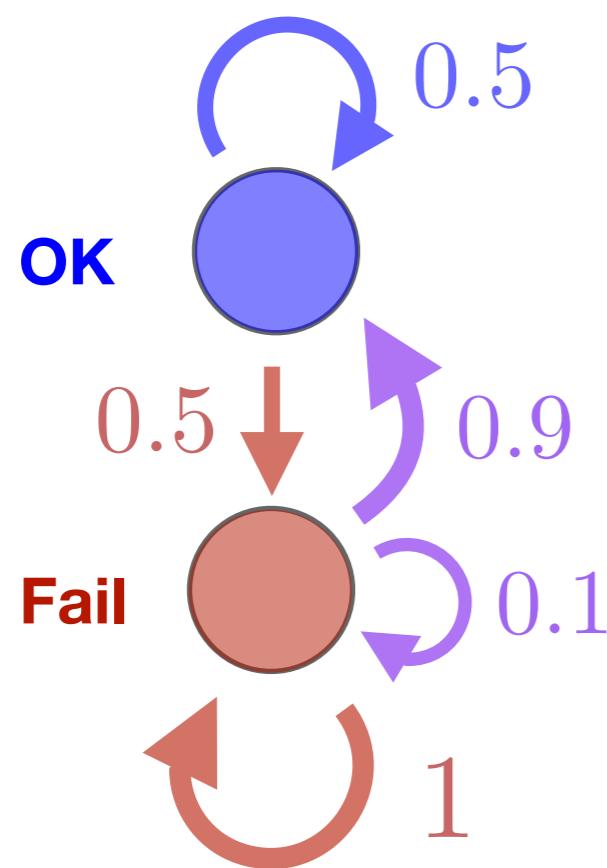
$t = 1$

$t = 2$

$t = 3$

$t = 4$

$t = 5$



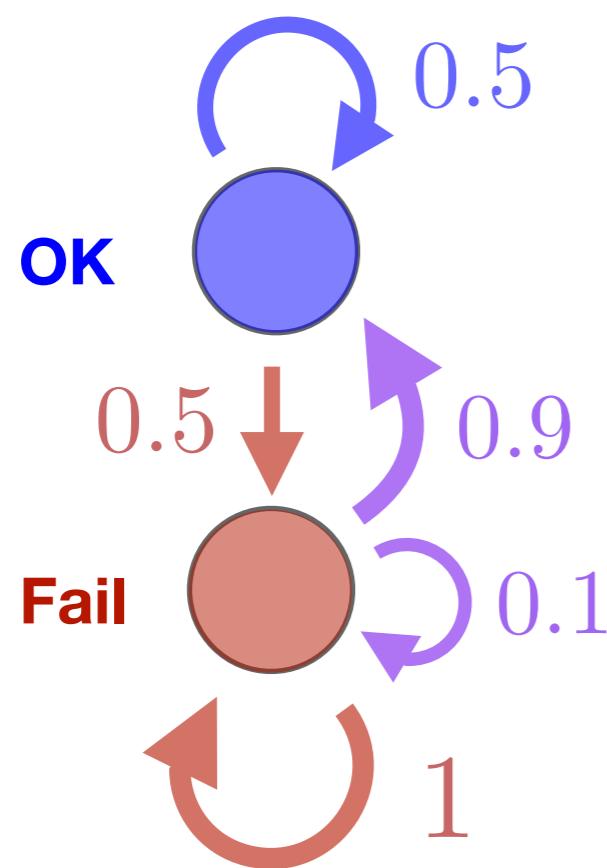
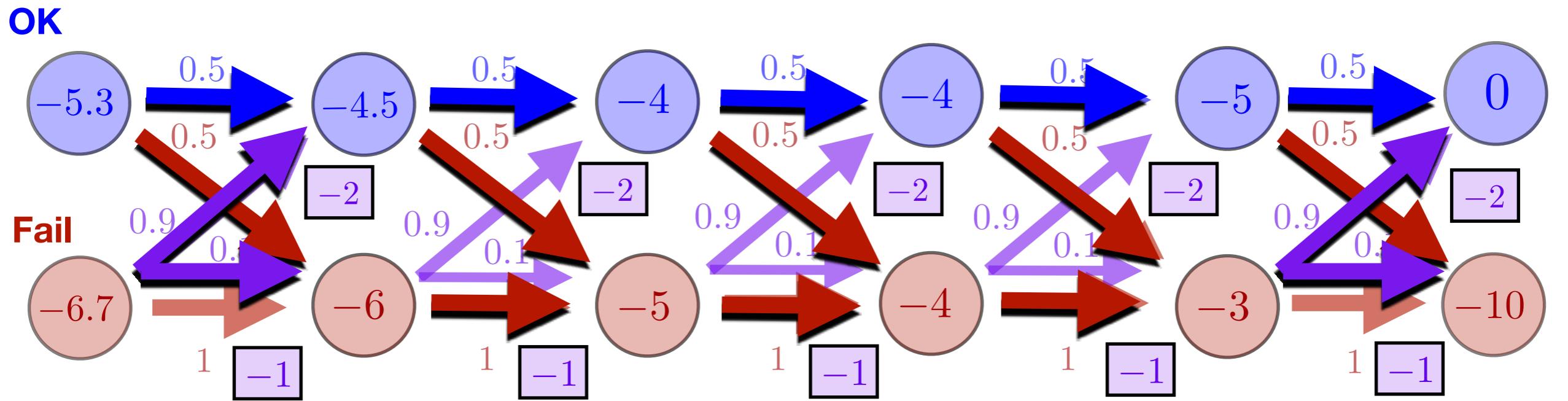
Discounted Bellman Equation

$$\gamma \in [0, 1)$$

$$v_s(\infty) = \max_{a \in \mathcal{A}_s} \sum_{s'} \gamma v_{s'}(\infty) P(s'|s, a) + r_{sa}$$

Discounted State-Value Meaning

$$v_s(\infty) = \frac{1}{1-\gamma} \sum_{k=0}^{\infty} \sum_{s'} \gamma^k r_{s'a} p_{s'a}(k), \quad p_s(0) = 1$$



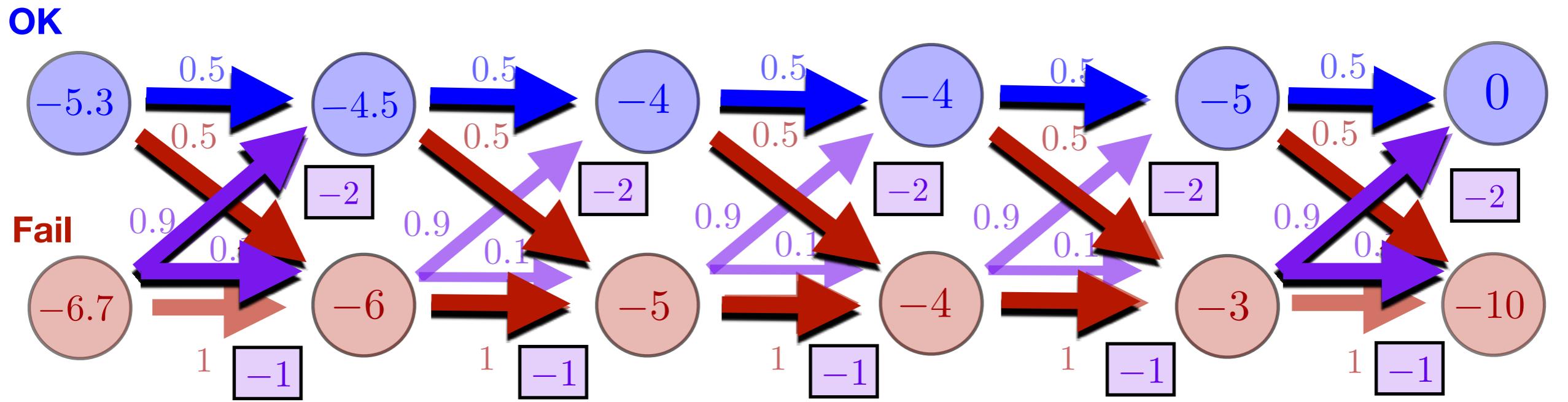
Discounted Bellman Equation

$$v_s(\infty) = \max_{a \in \mathcal{A}_s} \sum_{s'} \gamma v_{s'}(\infty) P(s'|s, a) + r_{sa}$$

Discounted Reward

$$R_\gamma = \sum_s p_s(0) v_s(\infty)$$

$$\gamma \in [0, 1)$$



$t = 0$

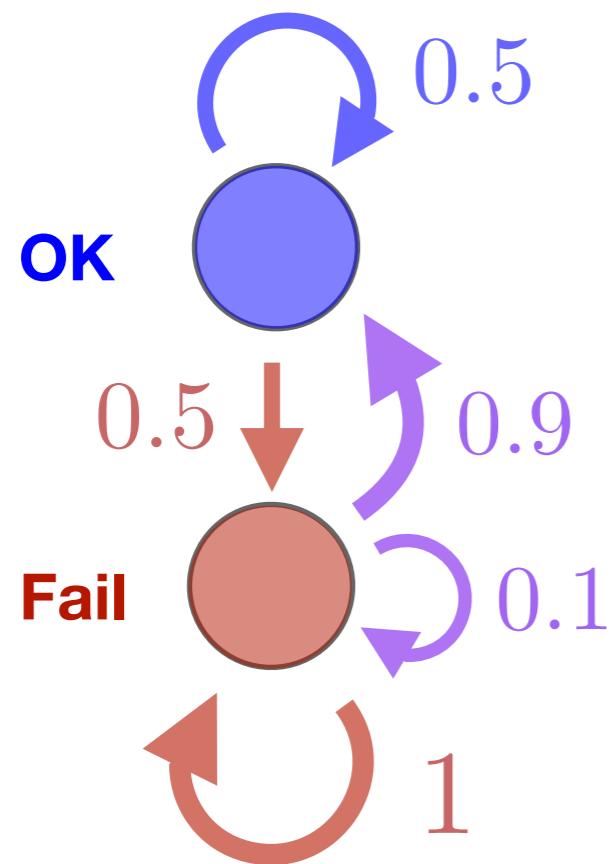
$t = 1$

$t = 2$

$t = 3$

$t = 4$

$t = 5$



Discounted Bellman Equation

$$v_s(\infty) = \max_{a \in \mathcal{A}_s} \sum_{s'} \gamma v_{s'}(\infty) P(s'|s, a) + r_{sa}$$

Discounted Reward

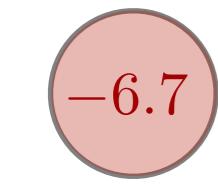
$$R_\gamma = \sum_s p_s(0) \left(\frac{1}{1-\gamma} \sum_{k=0}^{\infty} \sum_{s'} \gamma^k r_{s'a} p_{s'a}(k) \right)$$

$$\gamma \in [0, 1)$$

OK



Fail



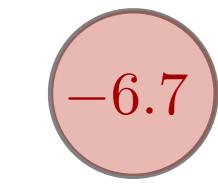
$t = 0$



$t = 1$



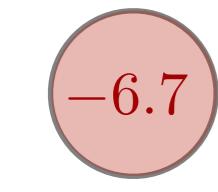
Fail



$t = 2$



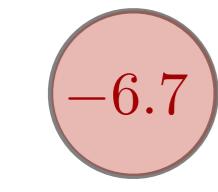
Fail



$t = 3$



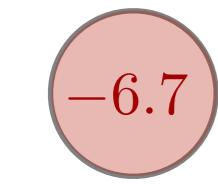
Fail



$t = 4$

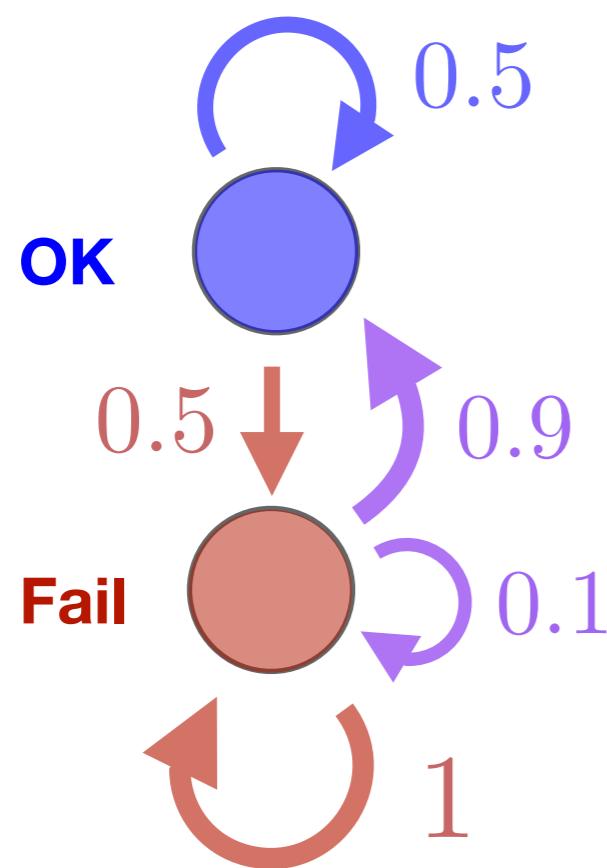


Fail



$t = 5$

Connection to Markov Chains



OK



Fail



$t = 0$



$t = 1$



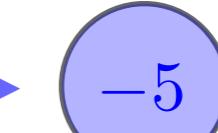
$t = 2$



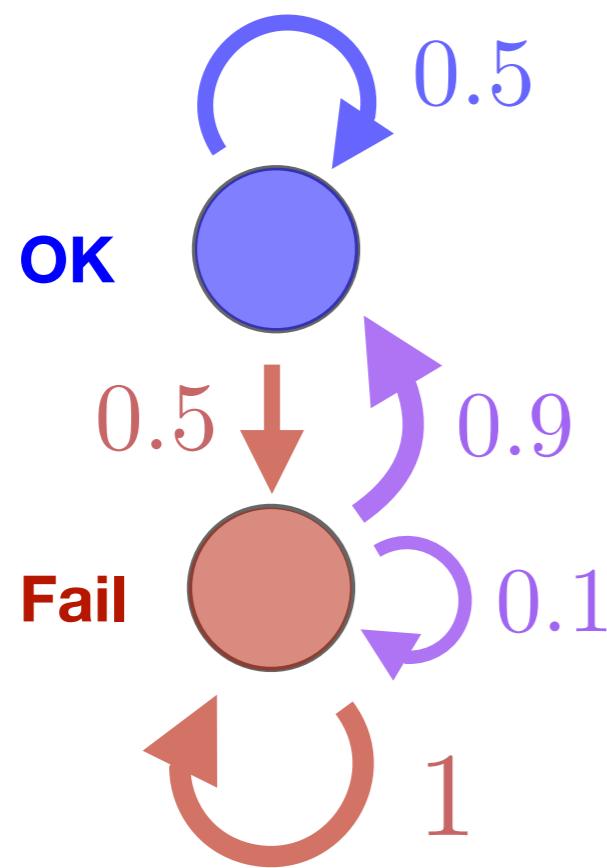
$t = 3$



$t = 4$



$t = 5$



Connection to Markov Chains

$$P = \left[\begin{array}{c|cc} \text{state 1} & \text{state 2} \\ \hline 0.5 & 0.9 & 0 \\ 0.5 & 0.1 & 1 \end{array} \right] \quad \begin{matrix} \text{action 1} & \text{action 1} & \text{action 2} \end{matrix}$$

*probability
transition
kernel*

OK



Fail



$t = 0$

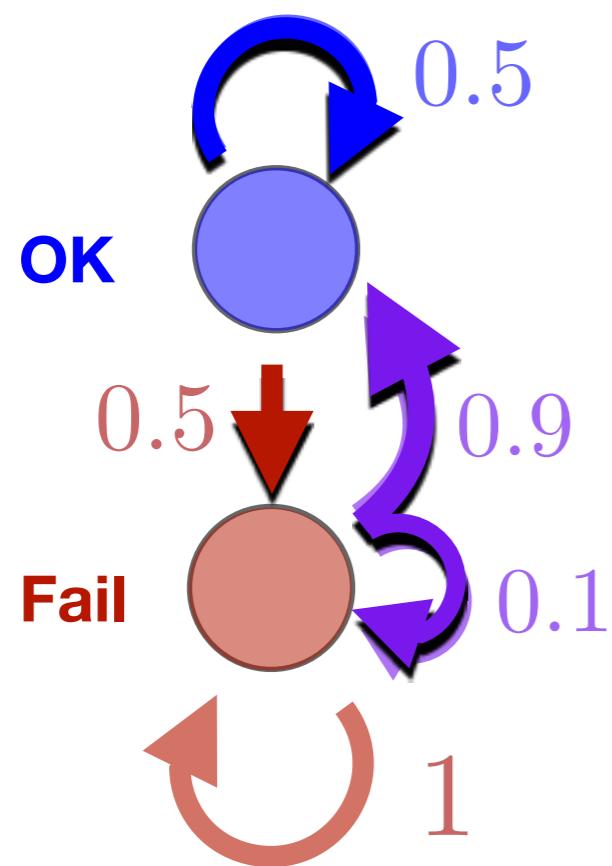
$t = 1$

$t = 2$

$t = 3$

$t = 4$

$t = 5$



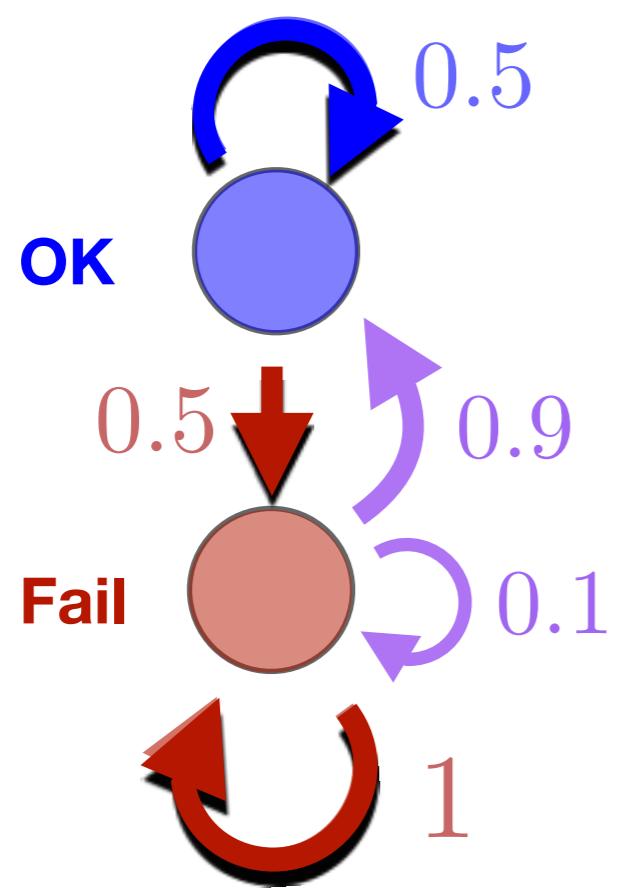
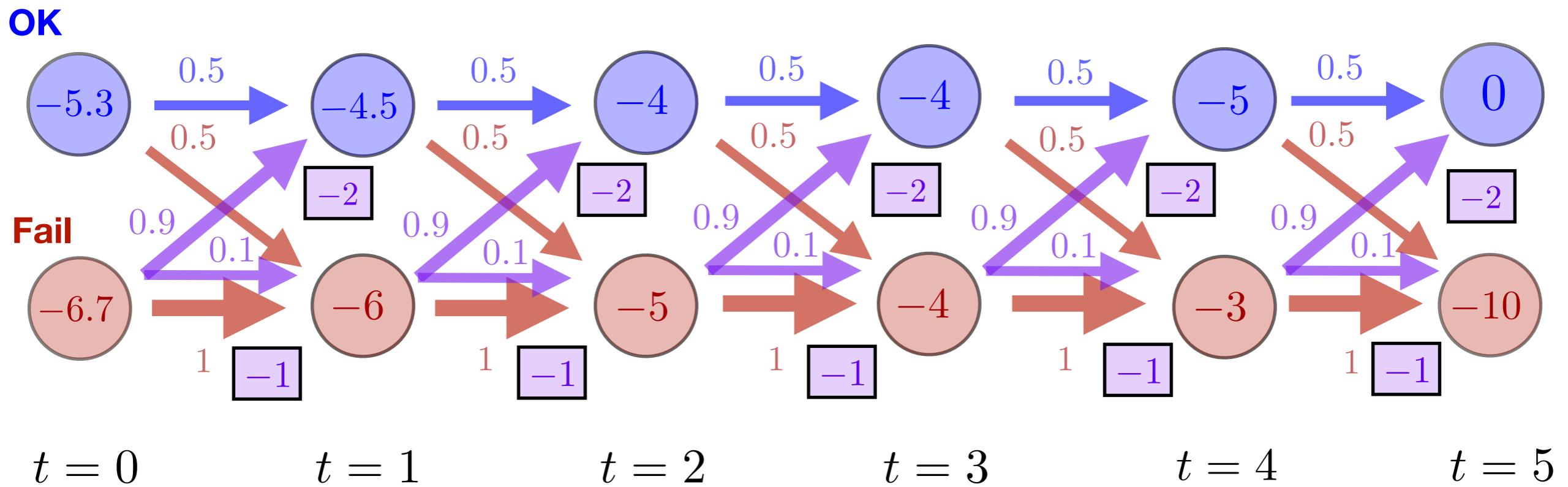
Connection to Markov Chains

$$P = \left[\begin{array}{c|cc} \text{state 1} & \text{state 2} \\ \hline 0.5 & 0.9 & 0 \\ 0.5 & 0.1 & 1 \end{array} \right]$$

probability transition kernel

$$M = \left[\begin{array}{c|c} 0.5 & 0.9 \\ 0.5 & 0.1 \end{array} \right]$$

“choosing a policy determines a Markov chain”



Connection to Markov Chains

$$P = \left[\begin{array}{c|cc} \text{state 1} & \text{state 2} \\ \hline 0.5 & 0.9 & 0 \\ 0.5 & 0.1 & 1 \end{array} \right]$$

action 1 action 1 action 2

$$M = \left[\begin{array}{c|c} 0.5 & 0 \\ 0.5 & 1 \end{array} \right]$$

*probability
transition
kernel*

*“choosing a
policy determines
a Markov chain”*