

6.036/6.862: Introduction to Machine Learning

Lecture: starts Tuesdays 9:35am (Boston time zone)

Course website: introml.odl.mit.edu

Who's talking? Prof. Tamara Broderick

Questions? discourse.odl.mit.edu ("Lecture 11" category)

Materials: Will all be available at course website

Last Time(s)

- I. State machines & Markov decision processes (MDPs)
- II. Choosing "best" actions
- 1 III. Value iteration; Q-learning

Today's Plan

- I. Back to supervised learning
- II. Sequential data
- III. Recurrent neural networks

Some terminology

Some terminology

- *Reinforcement learning* (RL): learning (to maximize rewards) by interacting with the world

Some terminology

- *Reinforcement learning* (RL): learning (to maximize rewards) by interacting with the world
- Contrast with *supervised learning*

Some terminology

- *Reinforcement learning* (RL): learning (to maximize rewards) by interacting with the world
 - Contrast with *supervised learning*
- *Model-based RL*: uses explicit conception of next state and reward given current state and action

Some terminology

- *Reinforcement learning* (RL): learning (to maximize rewards) by interacting with the world
 - Contrast with *supervised learning*
- *Model-based RL*: uses explicit conception of next state and reward given current state and action
 - “Model” used many different ways in machine learning

Some terminology

- *Reinforcement learning* (RL): learning (to maximize rewards) by interacting with the world
 - Contrast with *supervised learning*
- *Model-based RL*: uses explicit conception of next state and reward given current state and action
 - “Model” used many different ways in machine learning
 - Contrast with *Model-free RL*

Some terminology

- *Reinforcement learning* (RL): learning (to maximize rewards) by interacting with the world
 - Contrast with *supervised learning*
- *Model-based RL*: uses explicit conception of next state and reward given current state and action
 - “Model” used many different ways in machine learning
 - Contrast with *Model-free RL*
- *Q-learning*

Some terminology

- *Reinforcement learning* (RL): learning (to maximize rewards) by interacting with the world
 - Contrast with *supervised learning*
- *Model-based RL*: uses explicit conception of next state and reward given current state and action
 - “Model” used many different ways in machine learning
 - Contrast with *Model-free RL*
- *Q-learning*
 - Contrast with the Q^* function (expected reward of starting at s , making action a , and then making the “best” action ever after)

Some terminology

- *Reinforcement learning* (RL): learning (to maximize rewards) by interacting with the world
 - Contrast with *supervised learning*
- *Model-based RL*: uses explicit conception of next state and reward given current state and action
 - “Model” used many different ways in machine learning
 - Contrast with *Model-free RL*
- *Q-learning*
 - Contrast with the Q^* function (expected reward of starting at s , making action a , and then making the “best” action ever after)
 - Contrast with (any horizon) *value iteration*

Text prediction

Text prediction

Final product

VicePresidentOfCompany@HopefullyNotARealEmailA...

Final product

All the documents are finished. Please see attached

Text prediction

Final product

VicePresidentOfCompany@HopefullyNotARealEmailA..

Final product

All the documents are finished. Please see attached

WIKIPEDIA

The Free Encyclopedia

English
6 183 000+ articles

Español
1 637 000+ artículos

日本語
1 235 000+ 記事

Deutsch
2 495 000+ Artikel

Русский
1 672 000+ статей


Français
2 262 000+ articles

Italiano
1 645 000+ voci


中文
1 155 000+ 條目

Português
1 045 000+ artigos


Polski
1 435 000+ hasel




autocompEN



Autocomplete
Application that predicts the rest of a word a user is typing.




Search suggest drop-down list



Automotive industry in India

[ivoyage](#)
e travel guide



Automotive industry in the United States

[inews](#)
e news source

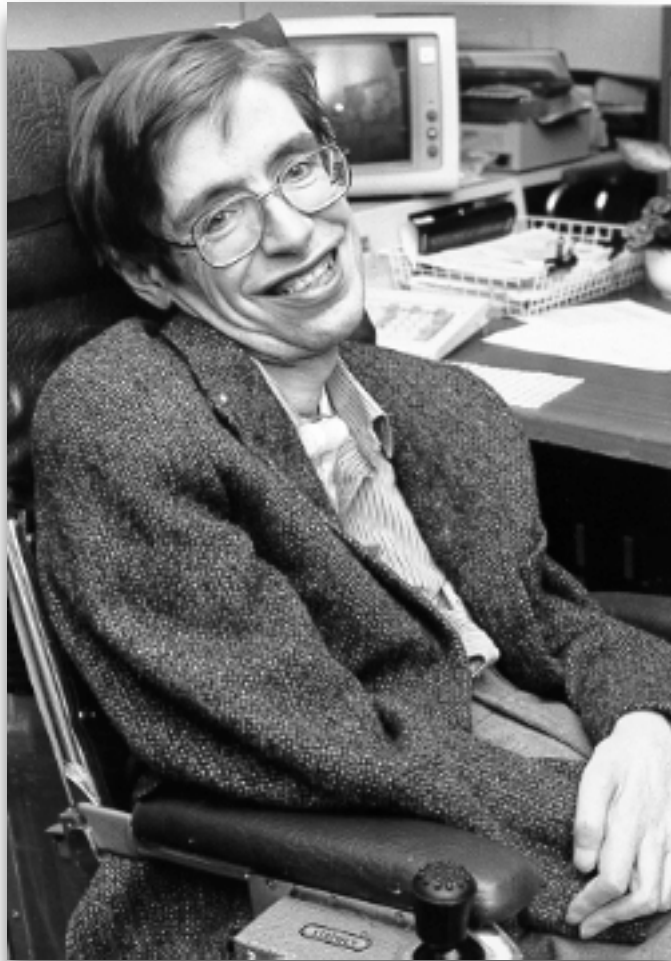
Text prediction

Final product

VicePresidentOfCompany@HopefullyNotARealEmailA..

Final product

All the documents are finished. Please see attached



WIKIPEDIA
The Free Encyclopedia

English

6 183 000+ articles

Español

1 637 000+ artículos

日本語

1 235 000+ 記事

Русский

1 672 000+ статей

Italiano

1 645 000+ voci



Deutsch

2 495 000+ Artikel

Français

2 262 000+ articles

中文

1 155 000+ 條目

Português

1 045 000+ artigos

Polski

1 435 000+ haseł

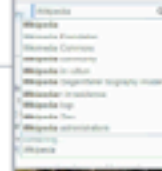
autocomp

EN



Autocomplete

Application that predicts the rest of a word a user is typing.



Search suggest drop-down list



Automotive industry in India



Automotive industry in the United States

ivoyage

e travel guide

inews

e news source

Text prediction

Final product

VicePresidentOfCompany@HopefullyNotARealEmailA..

Final product

All the documents are finished. Please see attached



WIKIPEDIA
The Free Encyclopedia

English

6 183 000+ articles

Español

1 637 000+ artículos

日本語

1 235 000+ 記事

Deutsch

2 495 000+ Artikel

Русский

1 672 000+ статей

Français

2 262 000+ articles

Italiano

1 645 000+ voci

中文

1 155 000+ 條目



Português

1 045 000+ artigos

Polski

1 435 000+ haseł

autocomp

EN



Autocomplete

Application that predicts the rest of a word a user is typing.

Search suggest drop-down list

Automotive industry in India



Automotive industry in the United States



ivoyage

e travel guide

inews

e news source

Text prediction: supervised learning

Text prediction: supervised learning

- Training data: lots of text

Text prediction: supervised learning

- Training data: lots of text
 - “what happens to a dream deferred”

Text prediction: supervised learning

- Training data: lots of text
 - “what happens to a dream deferred”

features	label
w	h

Text prediction: supervised learning

- Training data: lots of text
 - “what happens to a dream deferred”

features	label
w	h
wh	a

Text prediction: supervised learning

- Training data: lots of text
 - “what happens to a dream deferred”

features	label
w	h
wh	a
wha	t

Text prediction: supervised learning

- Training data: lots of text
 - “what happens to a dream deferred”

features	label
w	h
wh	a
wha	t
what	–
what_	h
what_h	a
what_ha	p
what_hap	p
what_happ	e

Text prediction: supervised learning

- Training data: lots of text
 - “what happens to a dream deferred”

features	label
w	h
wh	a
wha	t
what	–
what_	h
what_h	a
what_ha	p
what_hap	p
what_happ	e

- Classification with 27 classes

Text prediction: supervised learning

- Training data: lots of text
 - “what happens to a dream deferred”

features	label
w	h
wh	a
wha	t
what	–
what_	h
what_h	a
what_ha	p
what_hap	p
what_happ	e

- Classification with 27 classes
- How to featurize?

Text prediction: supervised learning

- Training data: lots of text
 - “what happens to a dream deferred”

features	label
w	h
wh	a
wha	t
what	–
what_	h
what_h	a
what_ha	p
what_hap	p
what_happ	e

- Classification with 27 classes
- How to featurize?
- Idea: use all previous characters.

Text prediction: supervised learning

- Training data: lots of text
 - “what happens to a dream deferred”

features	label
w	h
wh	a
wha	t
what	–
what_	h
what_h	a
what_ha	p
what_hap	p
what_happ	e

- Classification with 27 classes
- How to featurize?
- Idea: use all previous characters. But so far we’ve said $x^{(i)} \in \mathbb{R}^d$; i.e. fixed dimension

Text prediction: supervised learning

- Training data: lots of text
 - “what happens to a dream deferred”

features	label
w	h
wh	a
wha	t
what	–
what_	h
what_h	a
what_ha	p
what_hap	p
what_happ	e

- Classification with 27 classes
- How to featurize?
- Idea: use all previous characters. But so far we’ve said $x^{(i)} \in \mathbb{R}^d$; i.e. fixed dimension

Text prediction: supervised learning

- Training data: lots of text
 - “what happens to a dream deferred”

features	label
w	h
wh	a
wha	t
what	–
what_	h
what_h	a
what_ha	p
what_hap	p
what_happ	e

- Classification with 27 classes
- How to featurize?
- Idea: use all previous characters. But so far we’ve said $x^{(i)} \in \mathbb{R}^d$; i.e. fixed dimension

Text prediction: supervised learning

- Training data: lots of text
 - “what happens to a dream deferred”

features	label
w	h
wh	a
wha	t
what	–
what_	h
what_h	a
what_ha	p
what_hap	p
what_happ	e

- Classification with 27 classes
- How to featurize?
- Idea: use all previous characters. But so far we’ve said $x^{(i)} \in \mathbb{R}^d$; i.e. fixed dimension
- Idea: just use last character. But lose info

Text prediction: supervised learning

- Training data: lots of text
 - “what happens to a dream deferred”

features	label
w	h
wh	a
wha	t
what	–
what_	h
what_h	a
what_ha	p
what_hap	p
what_happ	e

- Classification with 27 classes
- How to featurize?
- Idea: use all previous characters. But so far we’ve said $x^{(i)} \in \mathbb{R}^d$; i.e. fixed dimension
- Idea: just use last character. But lose info
- Idea: use last m characters

Can express as a state machine

“wha”

Can express as a state machine

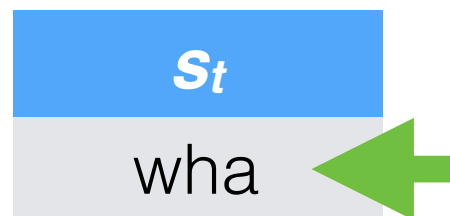
“wha”

Can express as a state machine

“wha”

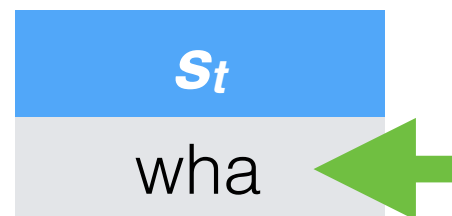
Can express as a state machine

“wha”



Can express as a state machine

“what”

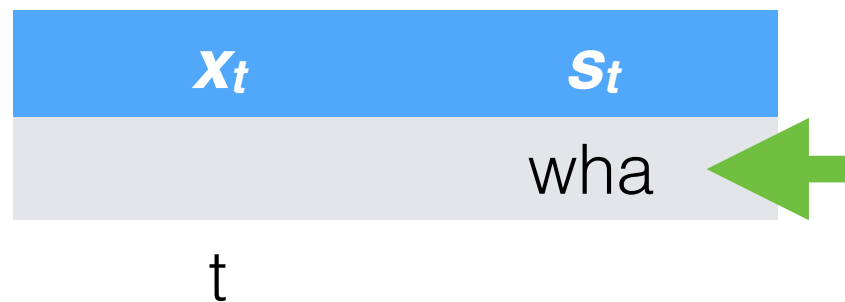


Can express as a state machine

“what”



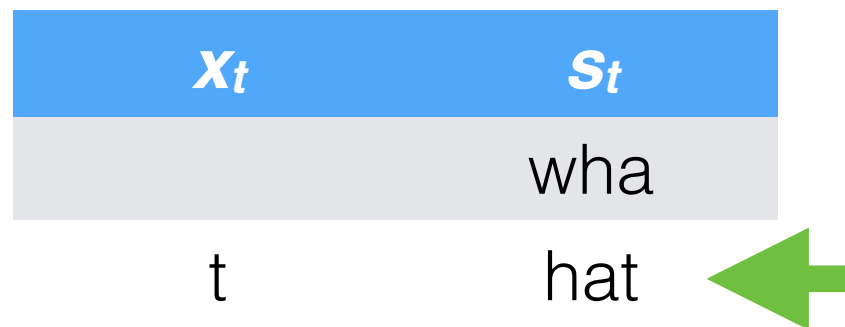
Can express as a state machine



“what”

Can express as a state machine


“what”



Can express as a state machine

“what_”


x_t	s_t
t	wha
	hat



Can express as a state machine

“what_”

x_t	s_t
t	wha
	hat



Can express as a state machine

“what_”

x_t	s_t
	wha
t	hat
_	at_



Can express as a state machine

“what happens to a
dream deferred”

<i>x_t</i>	<i>s_t</i>
	wha
t	hat
—	at_
h	t_h

Can express as a state machine

- Recall state machines:

“what happens to a
dream deferred”

x_t	s_t
	wha
t	hat
—	at_
h	t_h

Can express as a state machine

- Recall state machines:
 - Set of possible states \mathcal{S}

“what happens to a dream deferred”

x_t	s_t
	wha
t	hat
—	at_
h	t_h

Can express as a state machine

- Recall state machines:
 - Set of possible states \mathcal{S}
- Example:
 - All ordered m characters

“what happens to a
dream deferred”

x_t	s_t
	wha
t	hat
—	at_
h	t_h

Can express as a state machine

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
- Example:
 - All ordered m characters

“what happens to a
dream deferred”

x_t	s_t
	wha
t	hat
—	at_
h	t_h

Can express as a state machine

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
- Example:
 - All ordered m characters
 - All characters

“what happens to a
dream deferred”

x_t	s_t
	wha
t	hat
—	at_
h	t_h

Can express as a state machine

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
- Example:
 - All ordered m characters
 - All characters

“what happens to a
dream deferred”

x_t	s_t
	wha
t	hat
—	at_
h	t_h

Can express as a state machine

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
- Example:
 - All ordered m characters
 - All characters

“^what happens to a
dream deferred”

x_t	s_t
	wha
t	hat
—	at_
h	t_h

Can express as a state machine

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
- Example:
 - All ordered m characters
 - All characters

“^what happens to a
dream deferred”

x_t	s_t
	wha
t	hat
—	at_
h	t_h

Can express as a state machine

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
- Example:
 - All ordered m characters
 - All characters
 - m start characters

“^what happens to a dream deferred”

x_t	s_t
	wha
t	hat
—	at_
h	t_h

Can express as a state machine

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
- Example:
 - All ordered m characters
 - All characters
 - m start characters

x_t	s_t
	^^^
^	^^^
w	^^w
h	^wh
a	wha
t	hat
_	at_

“^what happens to a dream deferred”

Can express as a state machine

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
- Example:
 - All ordered m characters
 - All characters
 - m start characters

x_t	s_t
	^^^
^	^^^
w	^^w
h	^wh
a	wha
t	hat
_	at_

“^what happens to a dream deferred”

Can express as a state machine

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
- Example:
 - All ordered m characters
 - All characters
 - m start characters

x_t	s_t
	^^^
^	^^^
w	^^w
h	^wh
a	wha
t	hat
_	at_

“^what happens to a
dream deferred”



Can express as a state machine

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
- Example:
 - All ordered m characters
 - All characters
 - m start characters

x_t	s_t
	^^^
^	^^^
w	^^w
h	^wh
a	wha
t	hat
_	at_

“^what happens to a dream deferred”



Can express as a state machine

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
- Example:
 - All ordered m characters
 - All characters
 - m start characters

x_t	s_t
	^^^
^	^^^
w	^^w
h	^wh
a	wha
t	hat
_	at_

“^what happens to a dream deferred”

Can express as a state machine

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
- Example:
 - All ordered m characters
 - All characters
 - m start characters

x_t	s_t
	^^^
^	^^^
w	^^w
h	^wh
a	wha
t	hat
_	at_

“^what happens to a dream deferred”

Can express as a state machine

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
- Example:
 - All ordered m characters
 - All characters
 - m start characters

t	x_t	s_t
0		^^^
1	^	^^^
2	w	^^w
3	h	^wh
4	a	wha
5	t	hat
6	_	at_

“^what happens to a dream deferred”

Can express as a state machine

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
- Example:
 - All ordered m characters
 - All characters
 - m start characters

t	x_t	s_t
0		^^^
1	^	^^^
2	w	^^w
3	h	^wh
4	a	wha
5	t	hat
6	_	at_

“^what happens to a dream deferred”

Can express as a state machine

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
- Example:
 - All ordered m characters
 - All characters
 - m start characters
 - Update last m chars

t	x_t	s_t
0		^^^
1	^	^^^
2	w	^^w
3	h	^wh
4	a	wha
5	t	hat
6	_	at_

“^what happens to a dream deferred”

Can express as a state machine

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
- Example:
 - All ordered m characters
 - All characters
 - m start characters
 - Update last m chars

t	x_t	s_t
0		^^^
1	^	^^^
2	w	^^w
3	h	^wh
4	a	wha
5	t	hat
6	_	at_

“^what happens to a dream deferred”

Can express as a state machine

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
- Example:
 - All ordered m characters
 - All characters
 - m start characters
 - Update last m chars
 - All vectors of char probs

t	x_t	s_t
0		^^^
1	^	^^^
2	w	^^w
3	h	^wh
4	a	wha
5	t	hat
6	_	at_

“^what happens to a dream deferred”

Can express as a state machine

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $g(s)$
- Example:
 - All ordered m characters
 - All characters
 - m start characters
 - Update last m chars
 - All vectors of char probs

t	x_t	s_t
0		^^^
1	^	^^^
2	w	^^w
3	h	^wh
4	a	wha
5	t	hat
6	_	at_

“^what happens to a dream deferred”

Can express as a state machine

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $g(s)$
- Example:
 - All ordered m characters
 - All characters
 - m start characters
 - Update last m chars
 - All vectors of char probs
 - Multi-class linear classifier

t	x_t	s_t
0		^^^
1	^	^^^
2	w	^^w
3	h	^wh
4	a	wha
5	t	hat
6	_	at_

“^what happens to a dream deferred”

Can express as a state machine

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $g(s)$
- Example:
 - All ordered m characters
 - All characters
 - m start characters
 - Update last m chars
 - All vectors of char probs
 - Multi-class linear classifier

t	x_t	s_t
0		^^^
1	^	^^^
2	w	^^w
3	h	^wh
4	a	wha
5	t	hat
6	_	at_

- $x^{(1)}$: “^what happens to a dream deferred”

Can express as a state machine

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $g(s)$
- Example:
 - All ordered m characters
 - All characters
 - m start characters
 - Update last m chars
 - All vectors of char probs
 - Multi-class linear classifier

t	x_t	s_t
0		^^^
1	^	^^^
2	w	^^w
3	h	^wh
4	a	wha
5	t	hat
6	_	at_

- $x^{(1)}$: “^what happens to a dream deferred”
- $x^{(2)}$: “^if you can keep your head when all about you”
- $x^{(3)}$: “^you may write me down in history”

Can express as a state machine

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $g(s)$
- Example:
 - All ordered m characters
 - All characters
 - m start characters
 - Update last m chars
 - All vectors of char probs
 - Multi-class linear classifier

t	$x^{(1)}_t$	s_t
0		^^^
1	^	^^^
2	w	^^w
3	h	^wh
4	a	wha
5	t	hat
6	_	at_

- $x^{(1)}$: “^what happens to a dream deferred”
- $x^{(2)}$: “^if you can keep your head when all about you”
- $x^{(3)}$: “^you may write me down in history”

Can express as a state machine

t	x_t	s_t
0		$\wedge\wedge\wedge$
1	\wedge	$\wedge\wedge\wedge$
2	w	$\wedge\wedge w$
3	h	$\wedge wh$
4	a	wha
5	t	hat
6	_	at_

- $x^{(1)}$: “ \wedge what happens to a dream deferred”
- $x^{(2)}$: “ \wedge if you can keep your head when all about you”
- $x^{(3)}$: “ \wedge you may write me down in history”

Can express as a state machine

s_0

t	x_t	s_t
0		$\wedge\wedge\wedge$
1	\wedge	$\wedge\wedge\wedge$
2	w	$\wedge\wedge w$
3	h	$\wedge wh$
4	a	wha
5	t	hat
6	_	at_

- $x^{(1)}$: “ \wedge what happens to a dream deferred”
- $x^{(2)}$: “ \wedge if you can keep your head when all about you”
- $x^{(3)}$: “ \wedge you may write me down in history”

Can express as a state machine

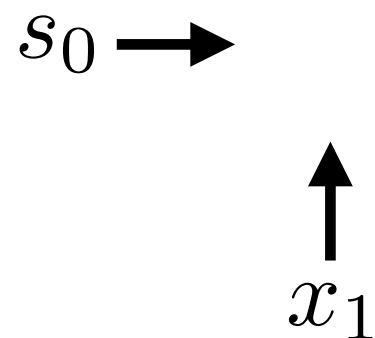
s_0

x_1

t	x_t	s_t
0		$\wedge \wedge \wedge$
1	\wedge	$\wedge \wedge \wedge$
2	w	$\wedge \wedge w$
3	h	$\wedge wh$
4	a	wha
5	t	hat
6	_	at_

- $x^{(1)}$: “ \wedge what happens to a dream deferred”
- $x^{(2)}$: “ \wedge if you can keep your head when all about you”
- $x^{(3)}$: “ \wedge you may write me down in history”

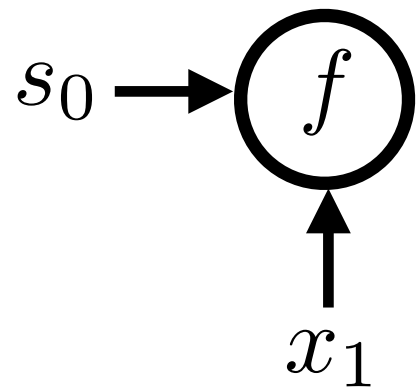
Can express as a state machine



t	x_t	s_t
0		^^^
1	^	^^^
2	w	^^w
3	h	^wh
4	a	wha
5	t	hat
6	_	at_

- $x^{(1)}$: “^what happens to a dream deferred”
- $x^{(2)}$: “^if you can keep your head when all about you”
- $x^{(3)}$: “^you may write me down in history”

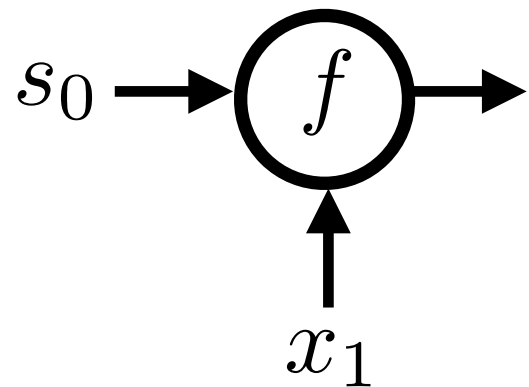
Can express as a state machine



t	x_t	s_t
0		$\wedge \wedge \wedge$
1	\wedge	$\wedge \wedge \wedge$
2	w	$\wedge \wedge w$
3	h	$\wedge wh$
4	a	wha
5	t	hat
6	_	at_

- $x^{(1)}$: “ \wedge what happens to a dream deferred”
- $x^{(2)}$: “ \wedge if you can keep your head when all about you”
- $x^{(3)}$: “ \wedge you may write me down in history”

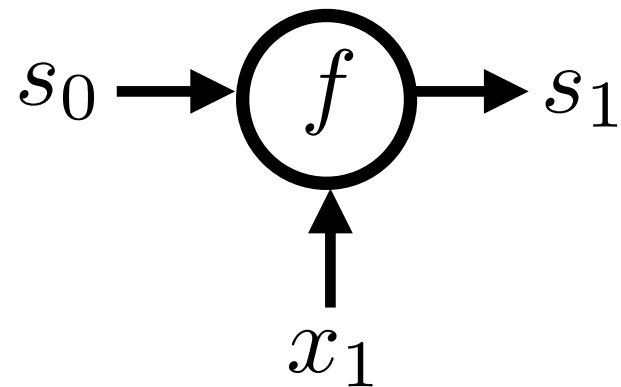
Can express as a state machine



t	x_t	s_t
0		$\wedge \wedge \wedge$
1	\wedge	$\wedge \wedge \wedge$
2	w	$\wedge \wedge w$
3	h	$\wedge wh$
4	a	wha
5	t	hat
6	_	at_

- $x^{(1)}$: “ \wedge what happens to a dream deferred”
- $x^{(2)}$: “ \wedge if you can keep your head when all about you”
- $x^{(3)}$: “ \wedge you may write me down in history”

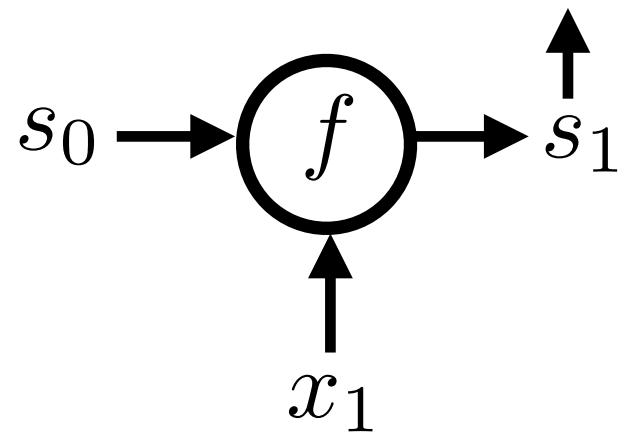
Can express as a state machine



t	x_t	s_t
0		^^^
1	^	^^^
2	w	^^w
3	h	^wh
4	a	wha
5	t	hat
6	_	at_

- $x^{(1)}$: “^what happens to a dream deferred”
- $x^{(2)}$: “^if you can keep your head when all about you”
- $x^{(3)}$: “^you may write me down in history”

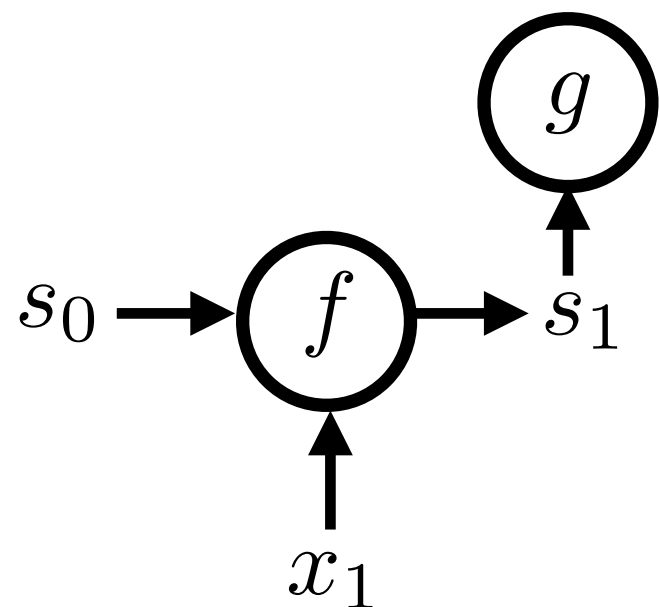
Can express as a state machine



t	x_t	s_t
0		^^^
1	^	^^^
2	w	^^w
3	h	^wh
4	a	wha
5	t	hat
6	_	at_

- $x^{(1)}$: “^what happens to a dream deferred”
- $x^{(2)}$: “^if you can keep your head when all about you”
- $x^{(3)}$: “^you may write me down in history”

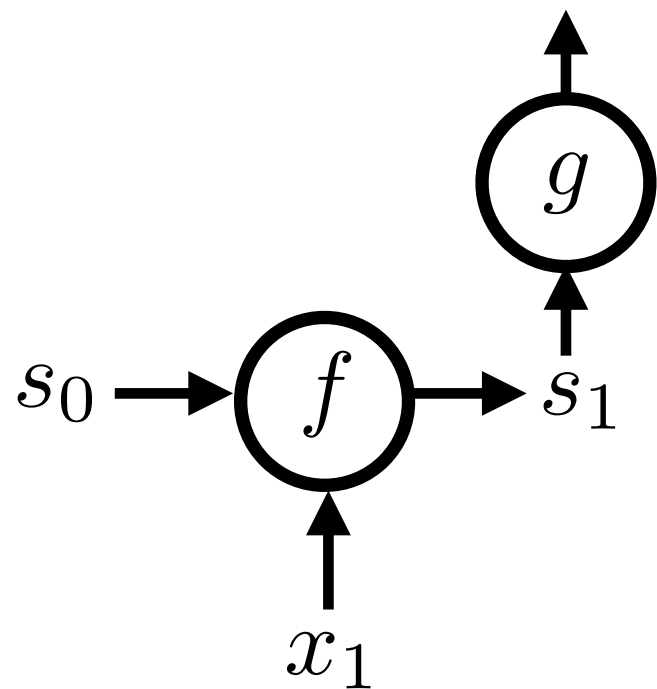
Can express as a state machine



t	x_t	s_t
0		$\wedge \wedge \wedge$
1	\wedge	$\wedge \wedge \wedge$
2	w	$\wedge \wedge w$
3	h	$\wedge wh$
4	a	wha
5	t	hat
6	_	at_

- $x^{(1)}$: “ \wedge what happens to a dream deferred”
- $x^{(2)}$: “ \wedge if you can keep your head when all about you”
- $x^{(3)}$: “ \wedge you may write me down in history”

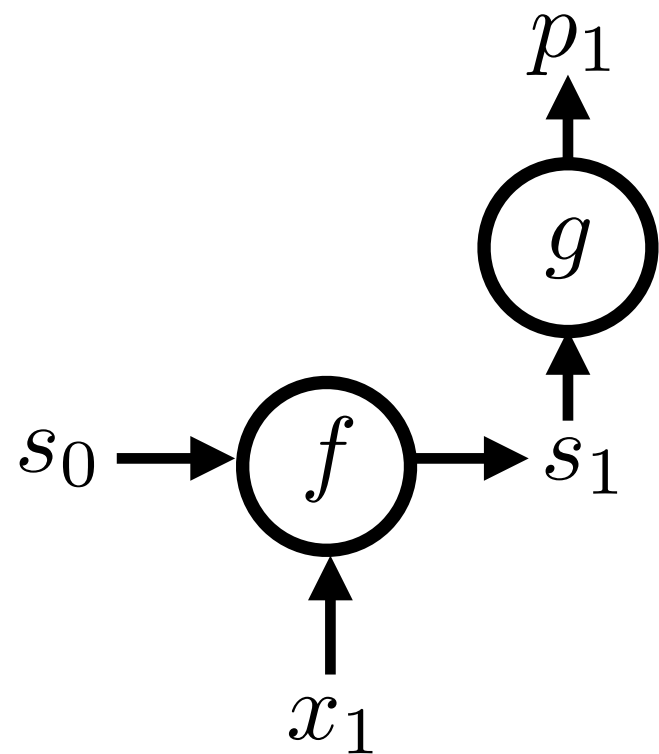
Can express as a state machine



t	x_t	s_t
0		$\wedge \wedge \wedge$
1	\wedge	$\wedge \wedge \wedge$
2	w	$\wedge \wedge w$
3	h	$\wedge wh$
4	a	wha
5	t	hat
6	_	at_

- $x^{(1)}$: “ \wedge what happens to a dream deferred”
- $x^{(2)}$: “ \wedge if you can keep your head when all about you”
- $x^{(3)}$: “ \wedge you may write me down in history”

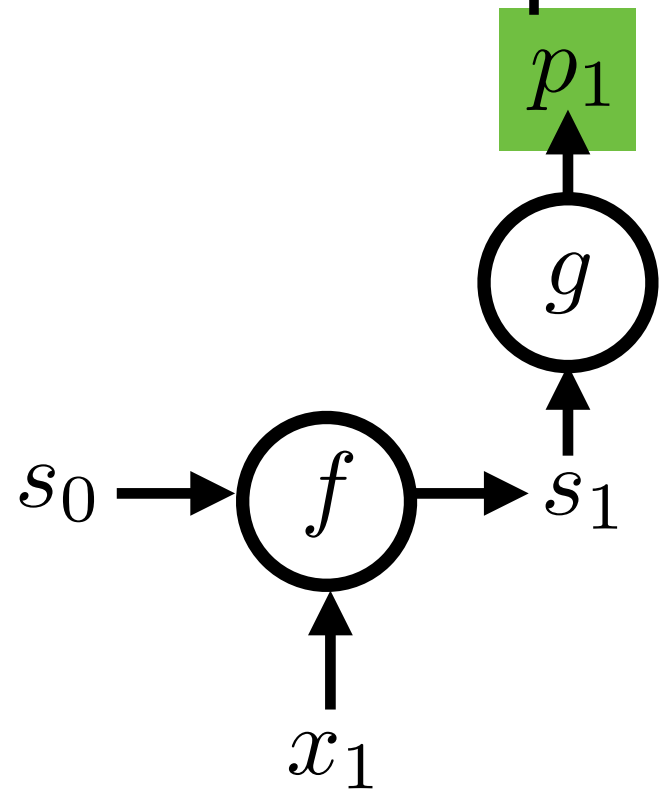
Can express as a state machine



t	x_t	s_t
0		$\wedge \wedge \wedge$
1	\wedge	$\wedge \wedge \wedge$
2	w	$\wedge \wedge w$
3	h	$\wedge wh$
4	a	wha
5	t	hat
6	_	at_

- $x^{(1)}$: “ \wedge what happens to a dream deferred”
- $x^{(2)}$: “ \wedge if you can keep your head when all about you”
- $x^{(3)}$: “ \wedge you may write me down in history”

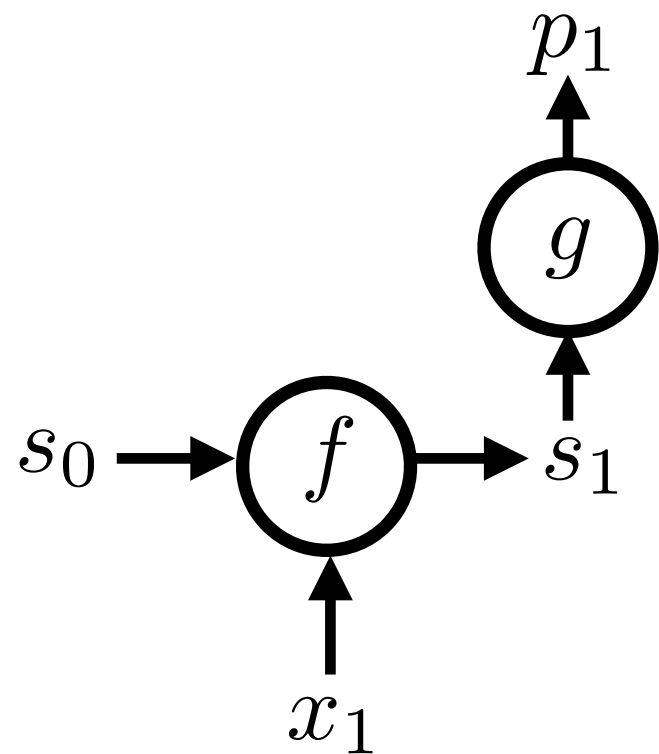
Can express as a state machine



t	x_t	s_t
0		$\wedge \wedge \wedge$
1	\wedge	$\wedge \wedge \wedge$
2	w	$\wedge \wedge w$
3	h	$\wedge wh$
4	a	wha
5	t	hat
6	_	at_

- $x^{(1)}$: “ \wedge what happens to a dream deferred”
- $x^{(2)}$: “ \wedge if you can keep your head when all about you”
- $x^{(3)}$: “ \wedge you may write me down in history”

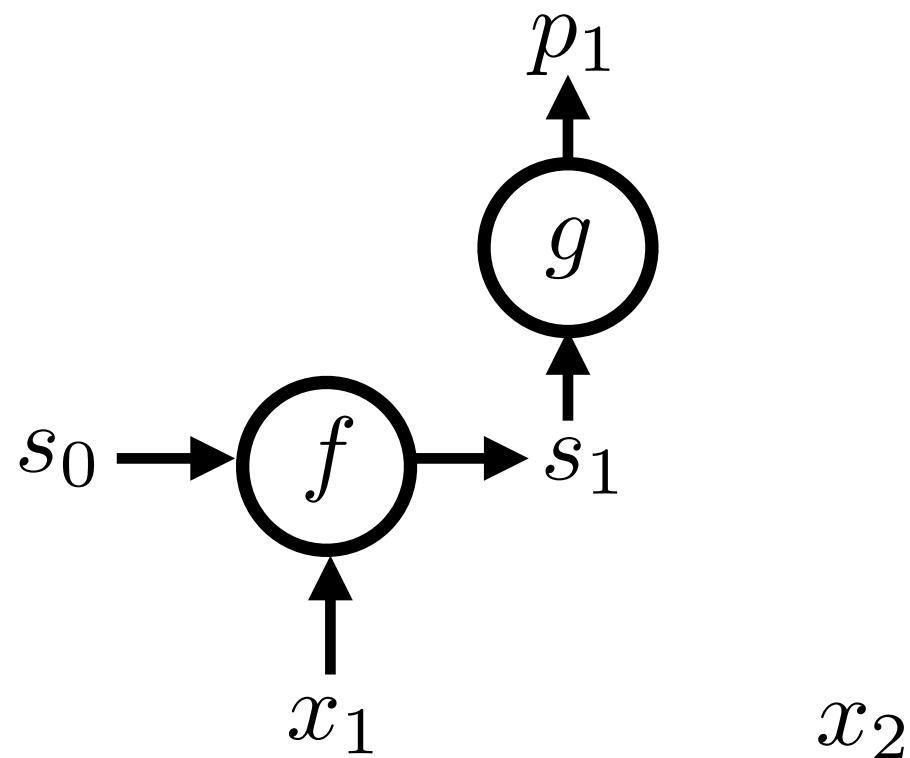
Can express as a state machine



t	x_t	s_t
0		$\wedge \wedge \wedge$
1	\wedge	$\wedge \wedge \wedge$
2	w	$\wedge \wedge w$
3	h	$\wedge wh$
4	a	wha
5	t	hat
6	_	at_

- $x^{(1)}$: “ \wedge what happens to a dream deferred”
- $x^{(2)}$: “ \wedge if you can keep your head when all about you”
- $x^{(3)}$: “ \wedge you may write me down in history”

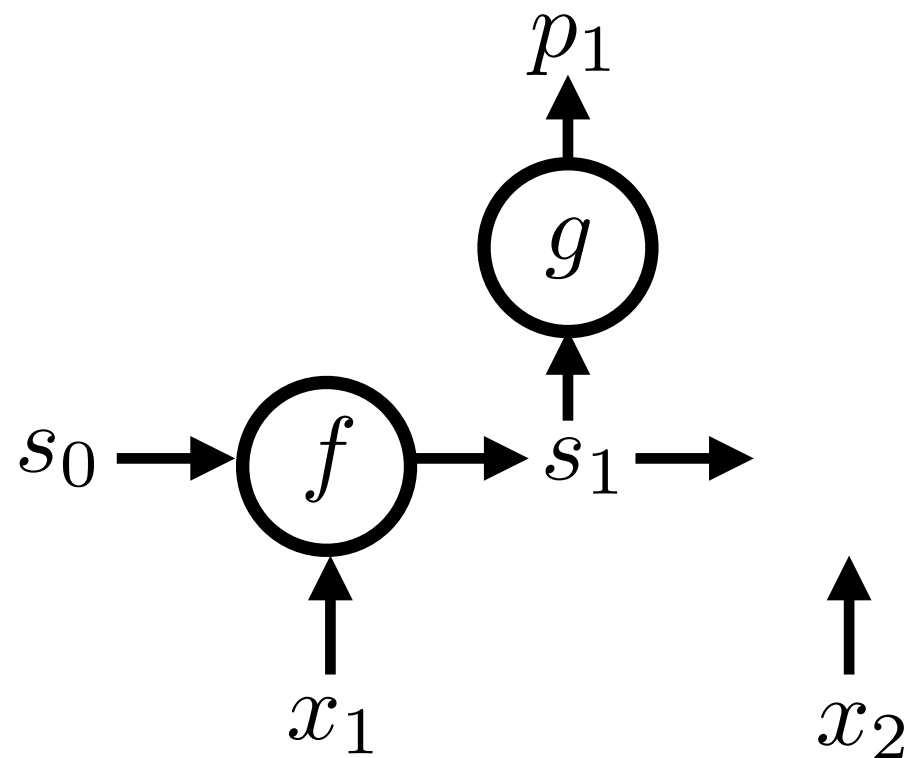
Can express as a state machine



t	x_t	s_t
0		$\wedge \wedge \wedge$
1	\wedge	$\wedge \wedge \wedge$
2	w	$\wedge \wedge w$
3	h	$\wedge wh$
4	a	wha
5	t	hat
6	_	at_

- $x^{(1)}$: “ \wedge what happens to a dream deferred”
- $x^{(2)}$: “ \wedge if you can keep your head when all about you”
- $x^{(3)}$: “ \wedge you may write me down in history”

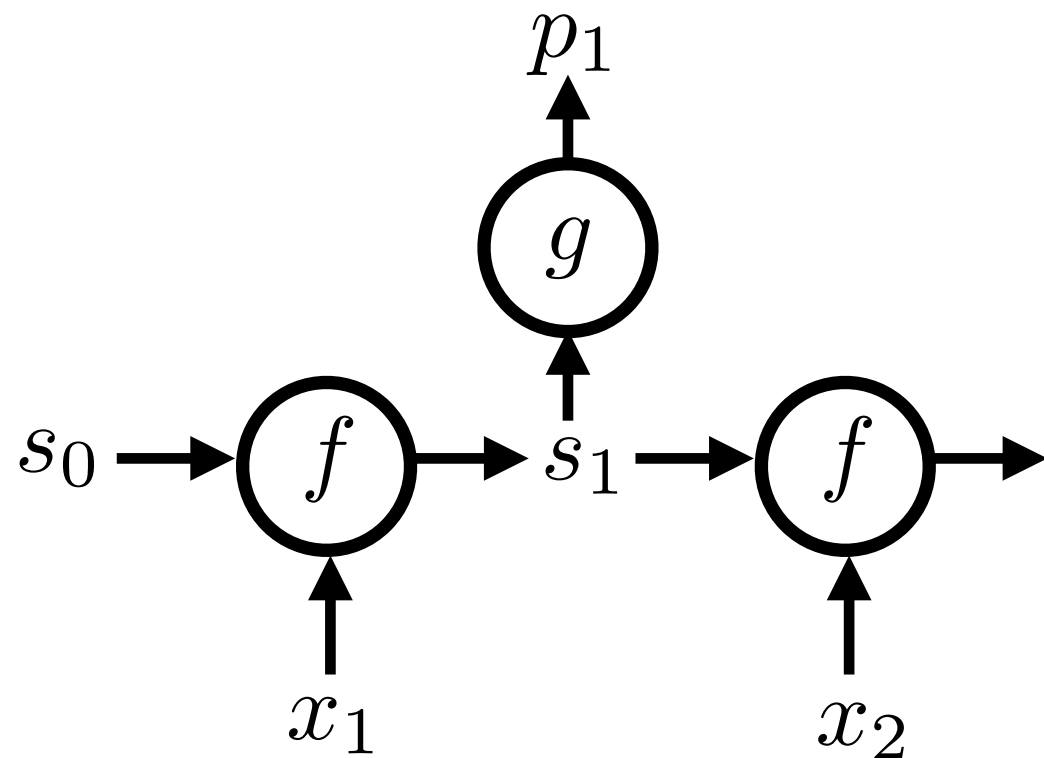
Can express as a state machine



t	x_t	s_t
0		^^^
1	^	^^^
2	w	^^w
3	h	^wh
4	a	wha
5	t	hat
6	_	at_

- $x^{(1)}$: “^what happens to a dream deferred”
- $x^{(2)}$: “^if you can keep your head when all about you”
- $x^{(3)}$: “^you may write me down in history”

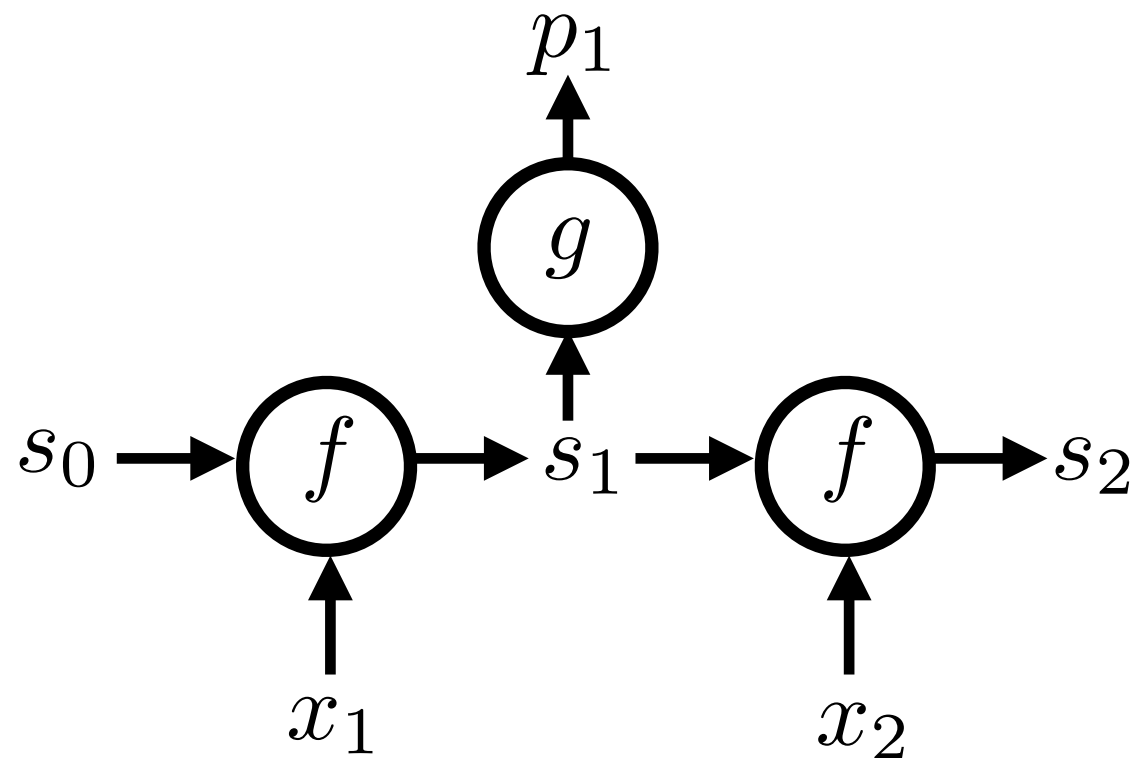
Can express as a state machine



t	x_t	s_t
0		$\wedge \wedge \wedge$
1	\wedge	$\wedge \wedge \wedge$
2	w	$\wedge \wedge w$
3	h	$\wedge wh$
4	a	wha
5	t	hat
6	_	at_

- $x^{(1)}$: “ \wedge what happens to a dream deferred”
- $x^{(2)}$: “ \wedge if you can keep your head when all about you”
- $x^{(3)}$: “ \wedge you may write me down in history”

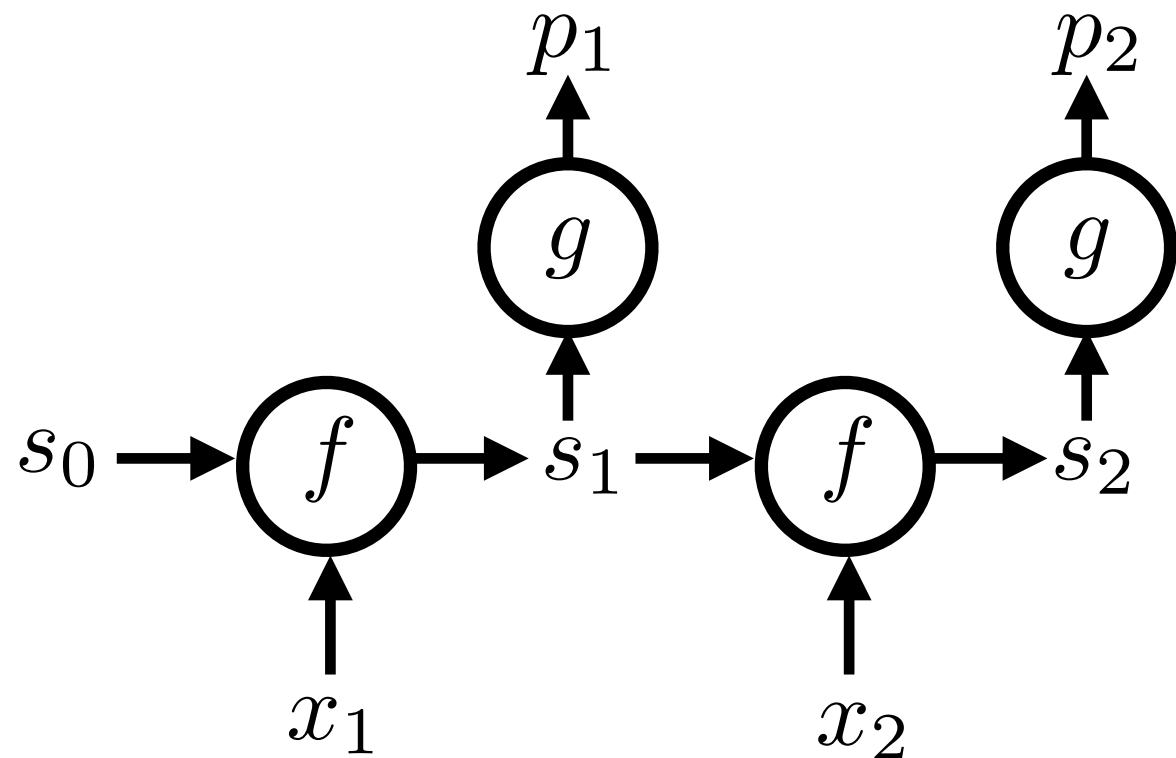
Can express as a state machine



t	x_t	s_t
0		^^^
1	^	^^^
2	w	^^w
3	h	^wh
4	a	wha
5	t	hat
6	_	at_

- $x^{(1)}$: “^what happens to a dream deferred”
- $x^{(2)}$: “^if you can keep your head when all about you”
- $x^{(3)}$: “^you may write me down in history”

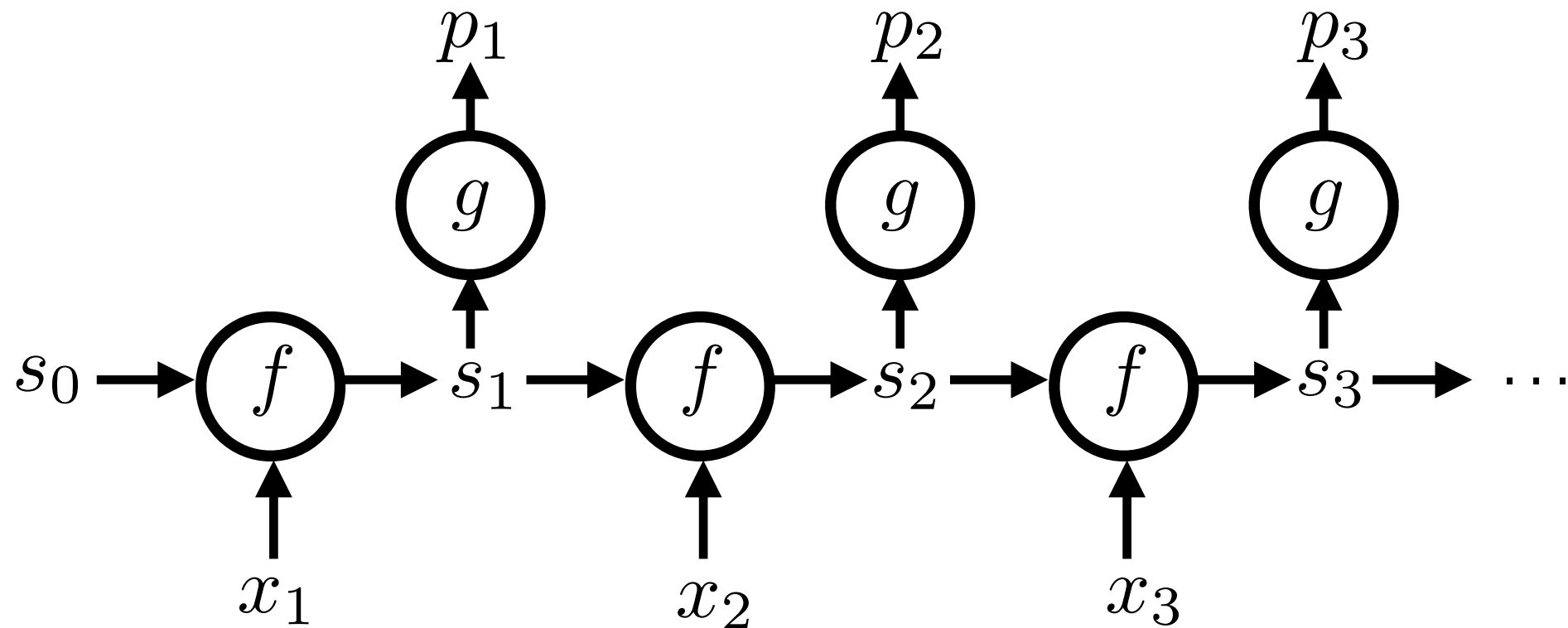
Can express as a state machine



t	x_t	s_t
0		^^^
1	^	^^^
2	w	^^w
3	h	^wh
4	a	wha
5	t	hat
6	_	at_

- $x^{(1)}$: “^what happens to a dream deferred”
- $x^{(2)}$: “^if you can keep your head when all about you”
- $x^{(3)}$: “^you may write me down in history”

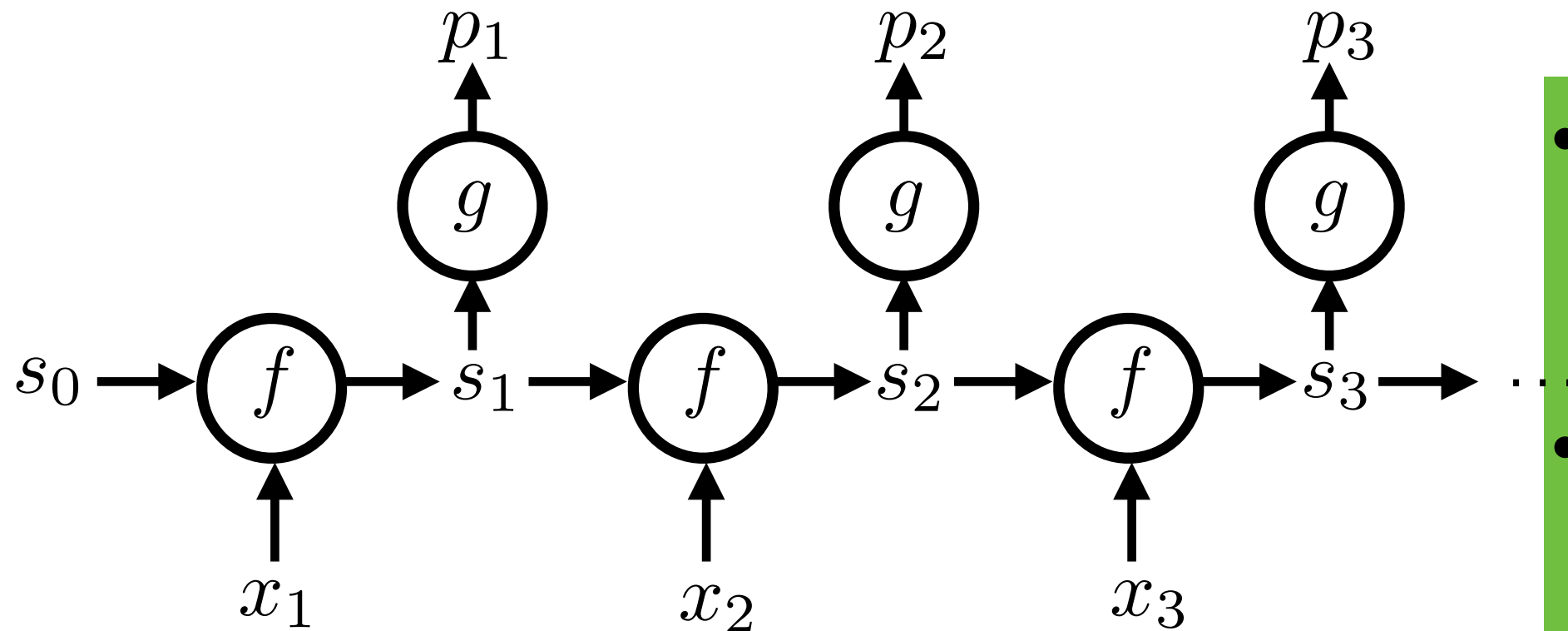
Can express as a state machine



t	x_t	s_t
0		^^^
1	^	^^^
2	w	^^w
3	h	^wh
4	a	wha
5	t	hat
6	_	at_

- $x^{(1)}$: “^what happens to a dream deferred”
- $x^{(2)}$: “^if you can keep your head when all about you”
- $x^{(3)}$: “^you may write me down in history”

Can express as a state machine

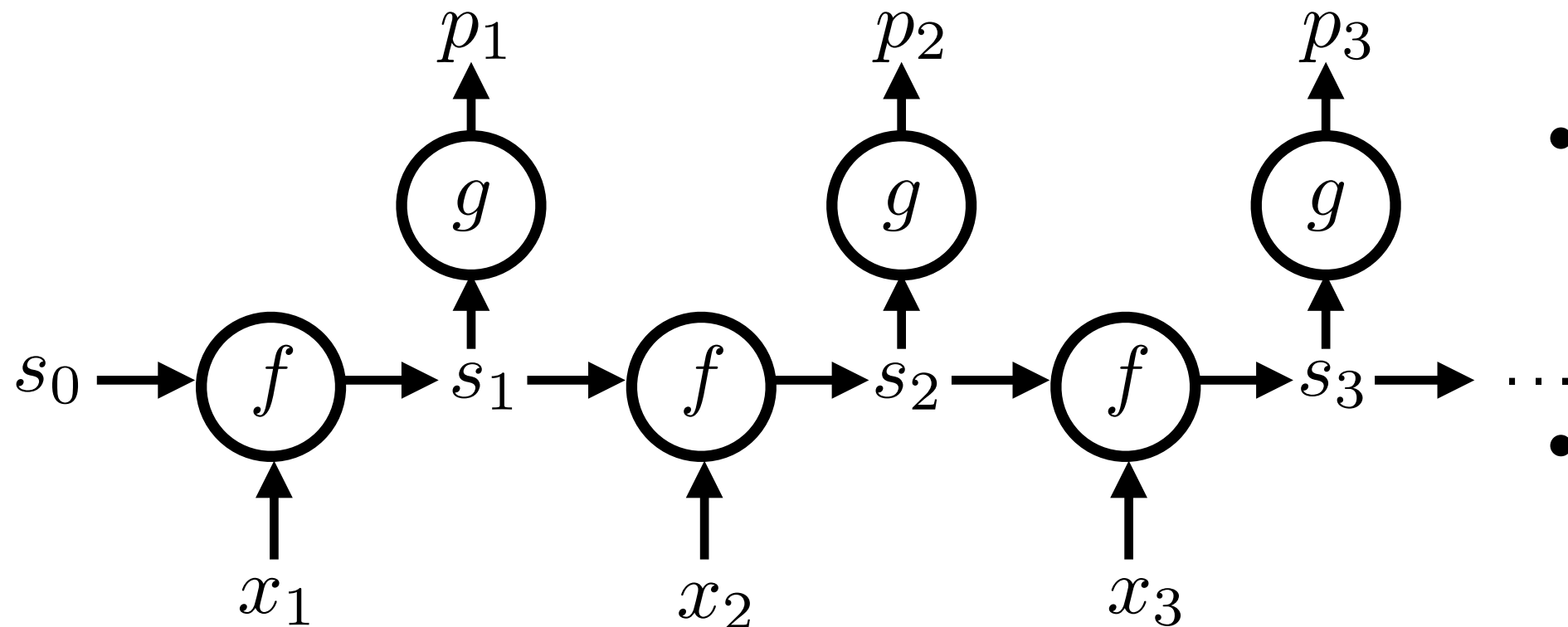


- m : number of characters in the context
- v : number of characters in the alphabet

t	x_t	s_t
0		^^^
1	^	^^^
2	w	^^w
3	h	^wh
4	a	wha
5	t	hat
6	_	at_

- $x^{(1)}$: “^what happens to a dream deferred”
- $x^{(2)}$: “^if you can keep your head when all about you”
- $x^{(3)}$: “^you may write me down in history”

Can express as a state machine

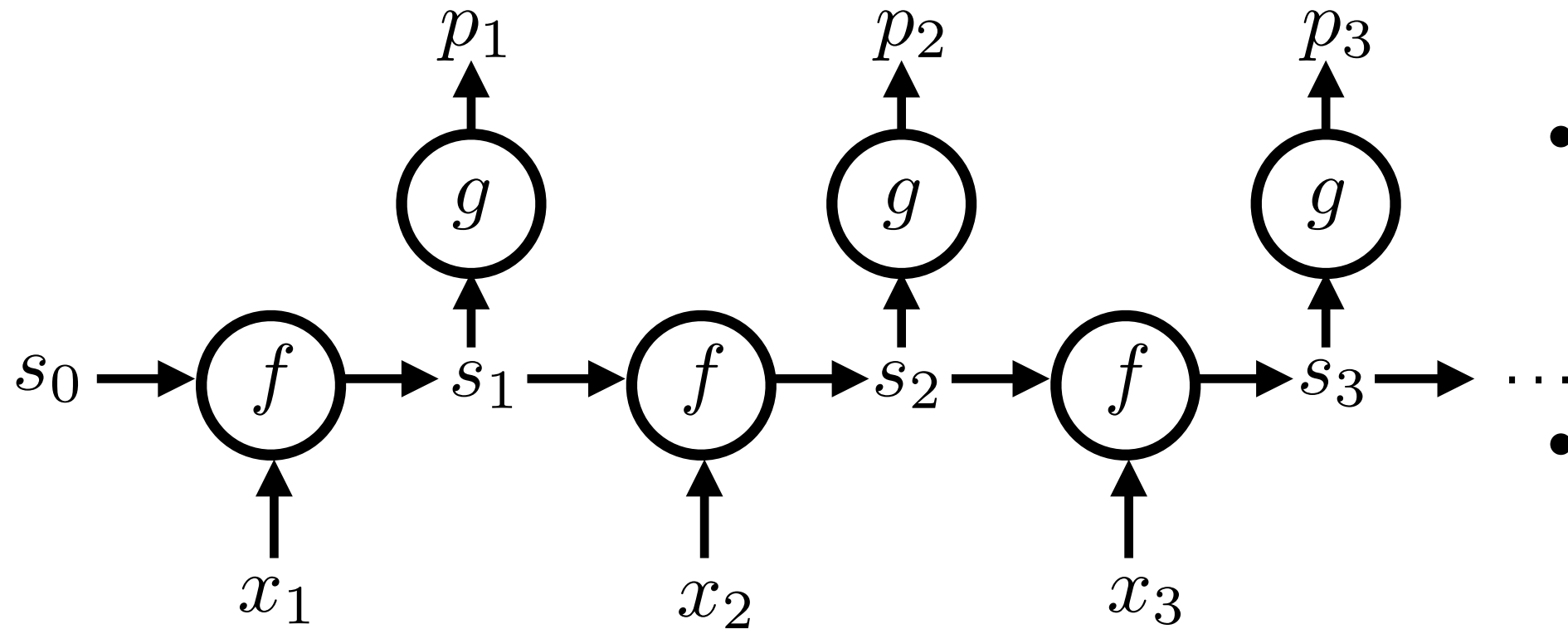


- m : number of characters in the context
- v : number of characters in the alphabet

t	x_t	s_t
0		^^^
1	^	^^^
2	w	^^w
3	h	^wh
4	a	wha
5	t	hat
6	_	at_

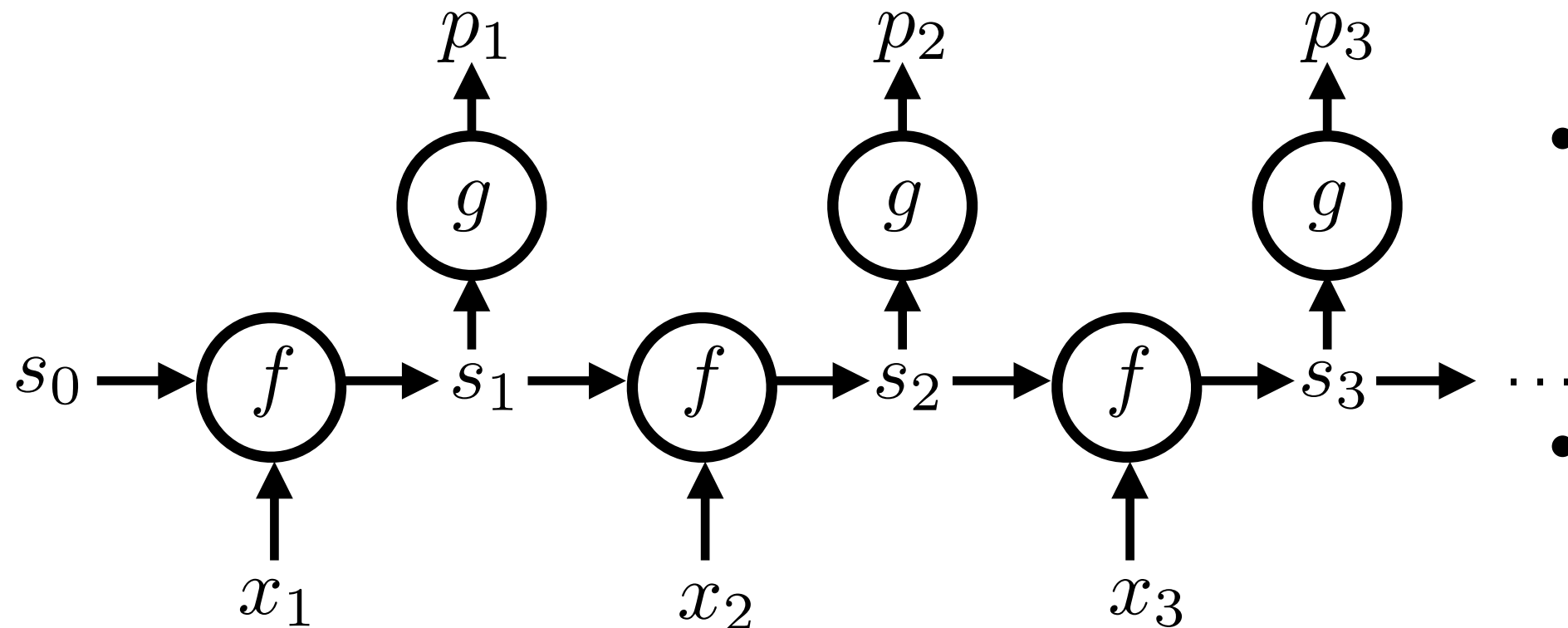
- $x^{(1)}$: “^what happens to a dream deferred”
- $x^{(2)}$: “^if you can keep your head when all about you”
- $x^{(3)}$: “^you may write me down in history”

Can express as a state machine



- m : number of characters in the context
- v : number of characters in the alphabet

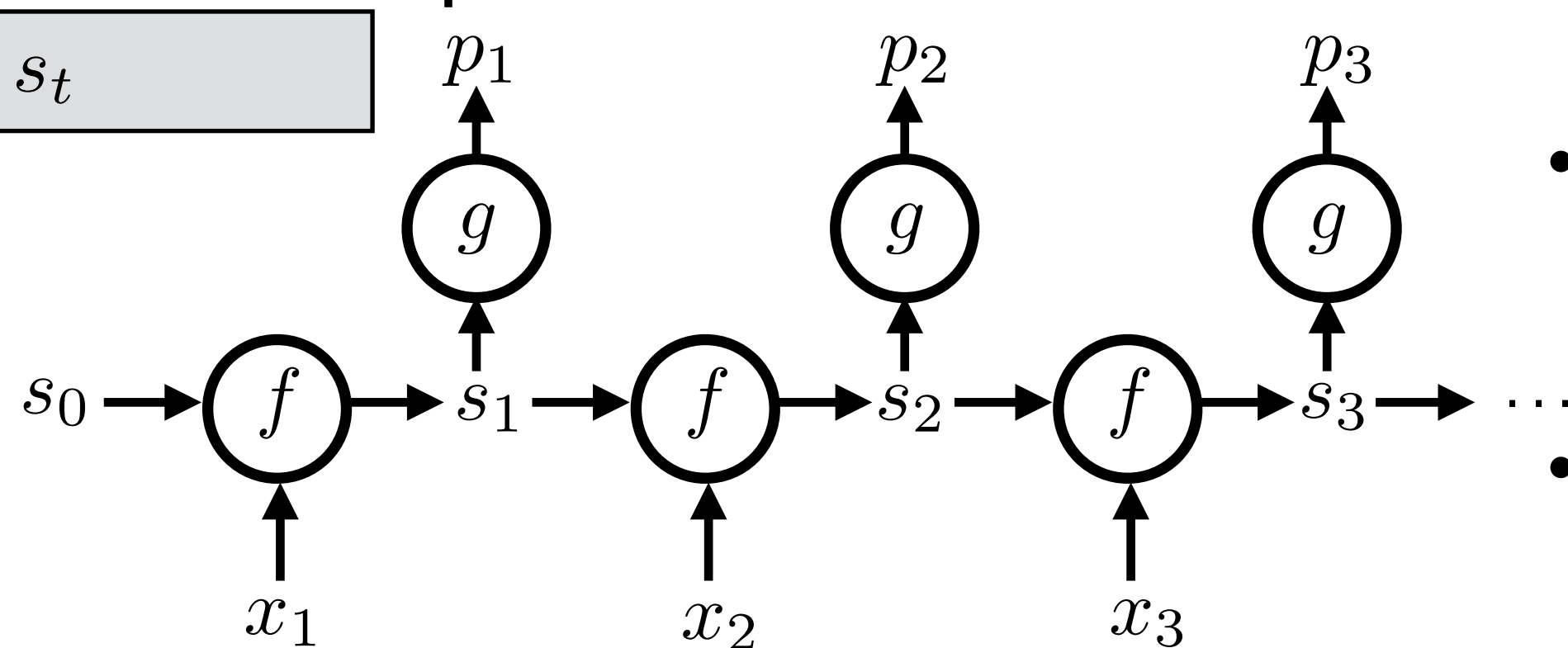
Can express as a state machine



- Example: Alphabet $\{0, 1\}$;
state is last $m = 3$ characters

- m : number of characters in the context
- v : number of characters in the alphabet

Can express as a state machine

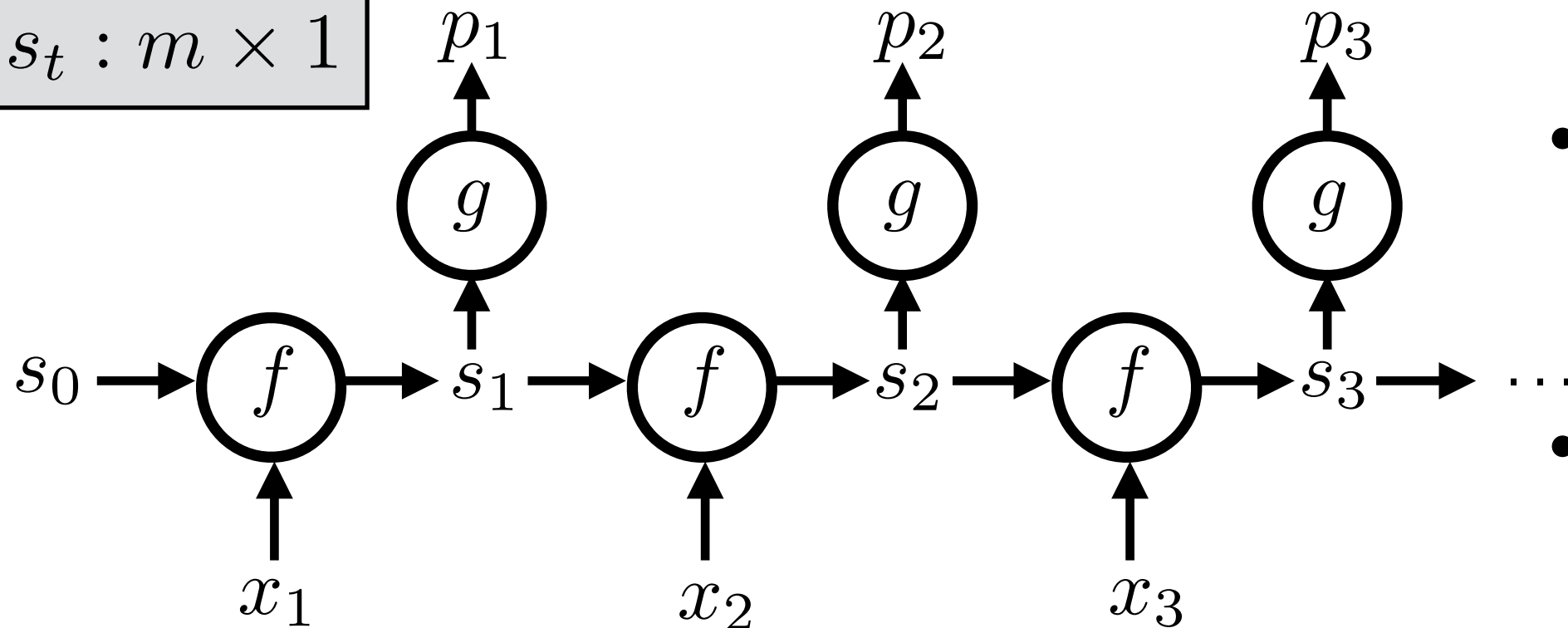


- Example: Alphabet $\{0, 1\}$;
state is last $m = 3$ characters

- m : number of characters in the context
- v : number of characters in the alphabet

Can express as a state machine

$$s_t : m \times 1$$

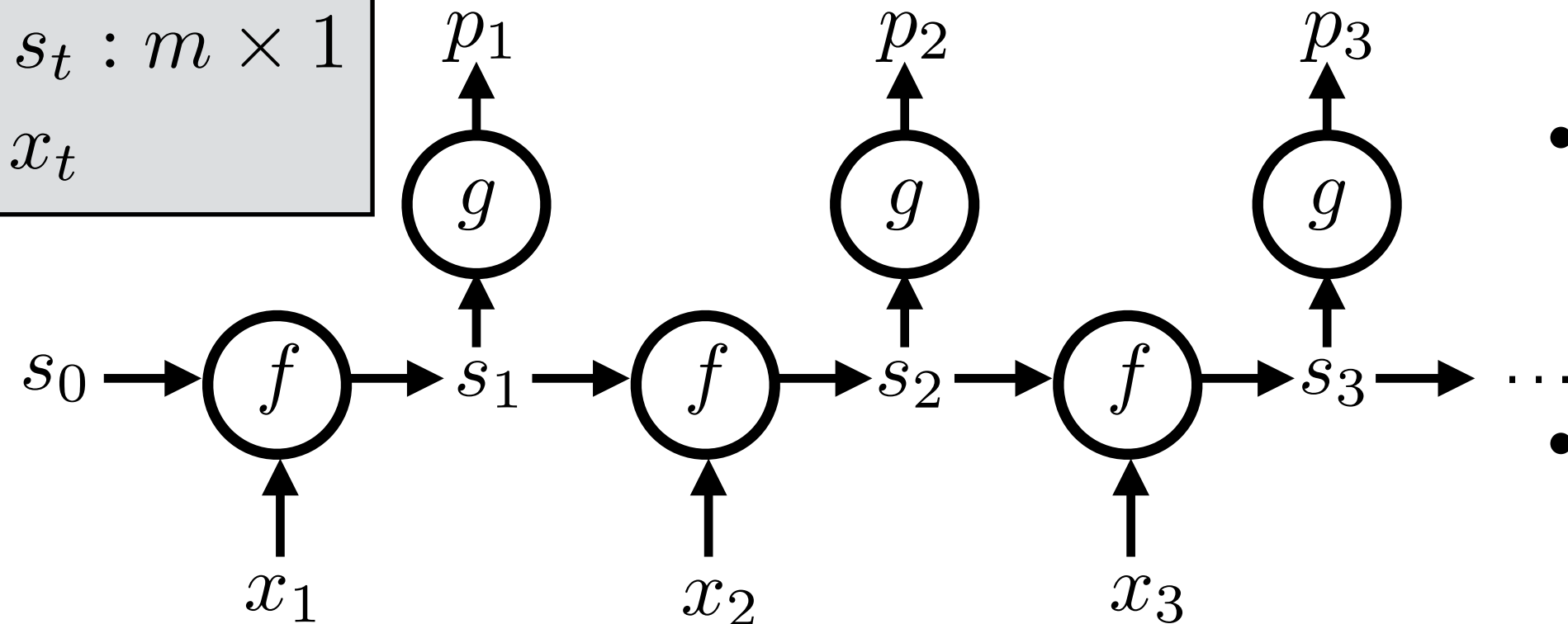


- Example: Alphabet $\{0, 1\}$;
state is last $m = 3$ characters

- m : number of characters in the context
- v : number of characters in the alphabet

Can express as a state machine

$$\begin{matrix} s_t : m \times 1 \\ x_t \end{matrix}$$

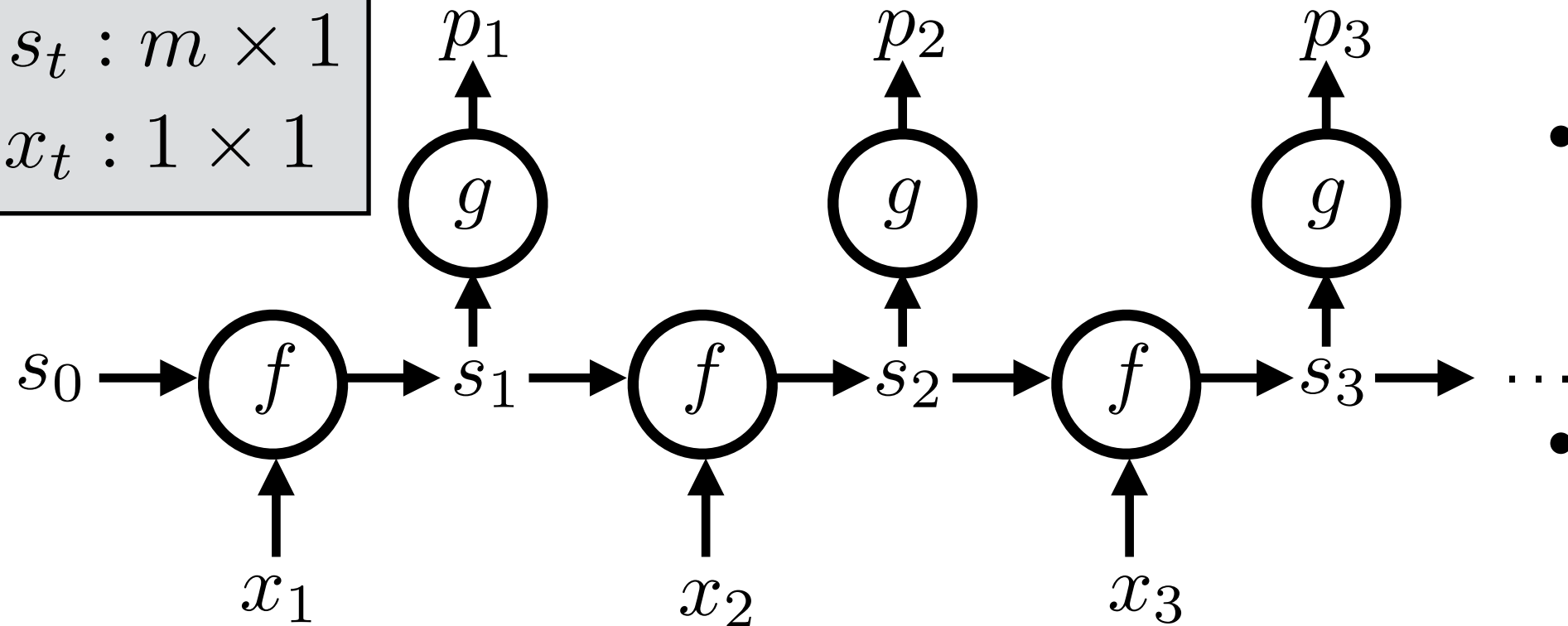


- Example: Alphabet $\{0, 1\}$; state is last $m = 3$ characters

- m : number of characters in the context
- v : number of characters in the alphabet

Can express as a state machine

$$\begin{matrix} s_t : m \times 1 \\ x_t : 1 \times 1 \end{matrix}$$

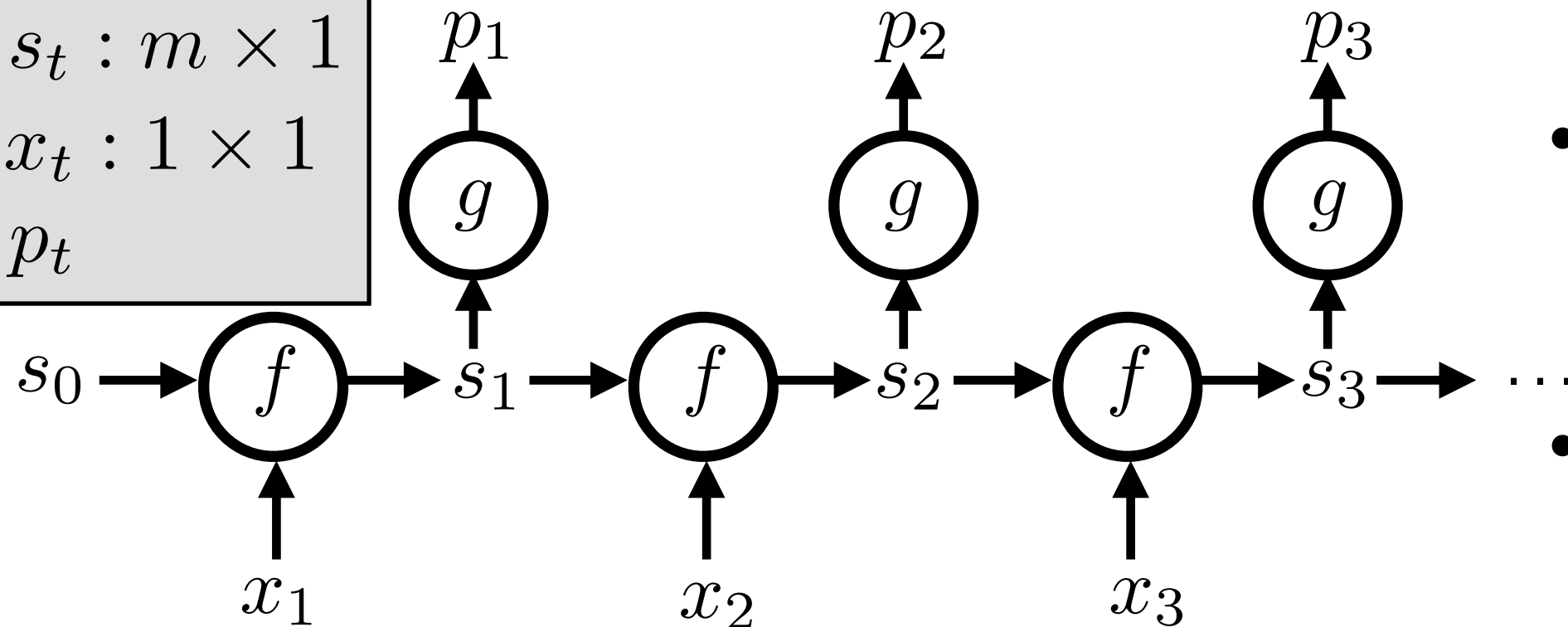


- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

- m : number of characters in the context
- v : number of characters in the alphabet

Can express as a state machine

$$\begin{array}{l} s_t : m \times 1 \\ x_t : 1 \times 1 \\ p_t \end{array}$$

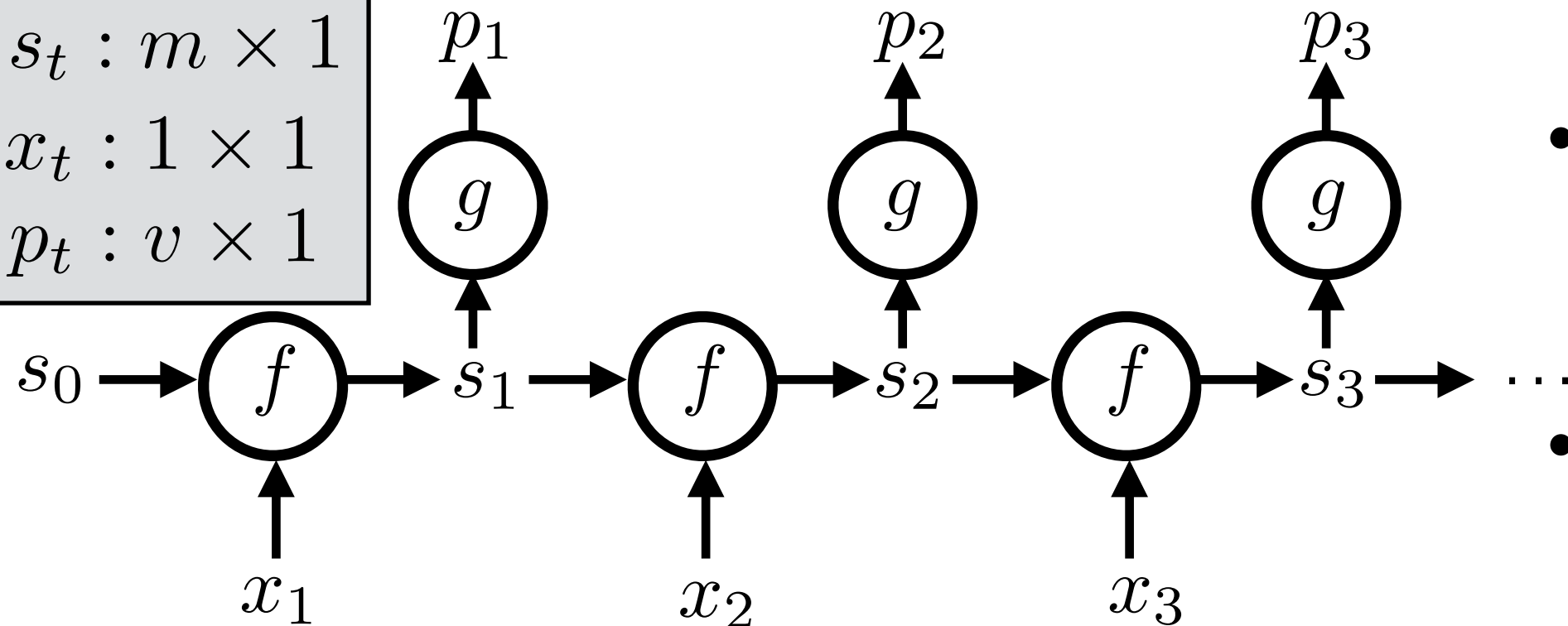


- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

- m : number of characters in the context
- v : number of characters in the alphabet

Can express as a state machine

$$\begin{array}{l} s_t : m \times 1 \\ x_t : 1 \times 1 \\ p_t : v \times 1 \end{array}$$

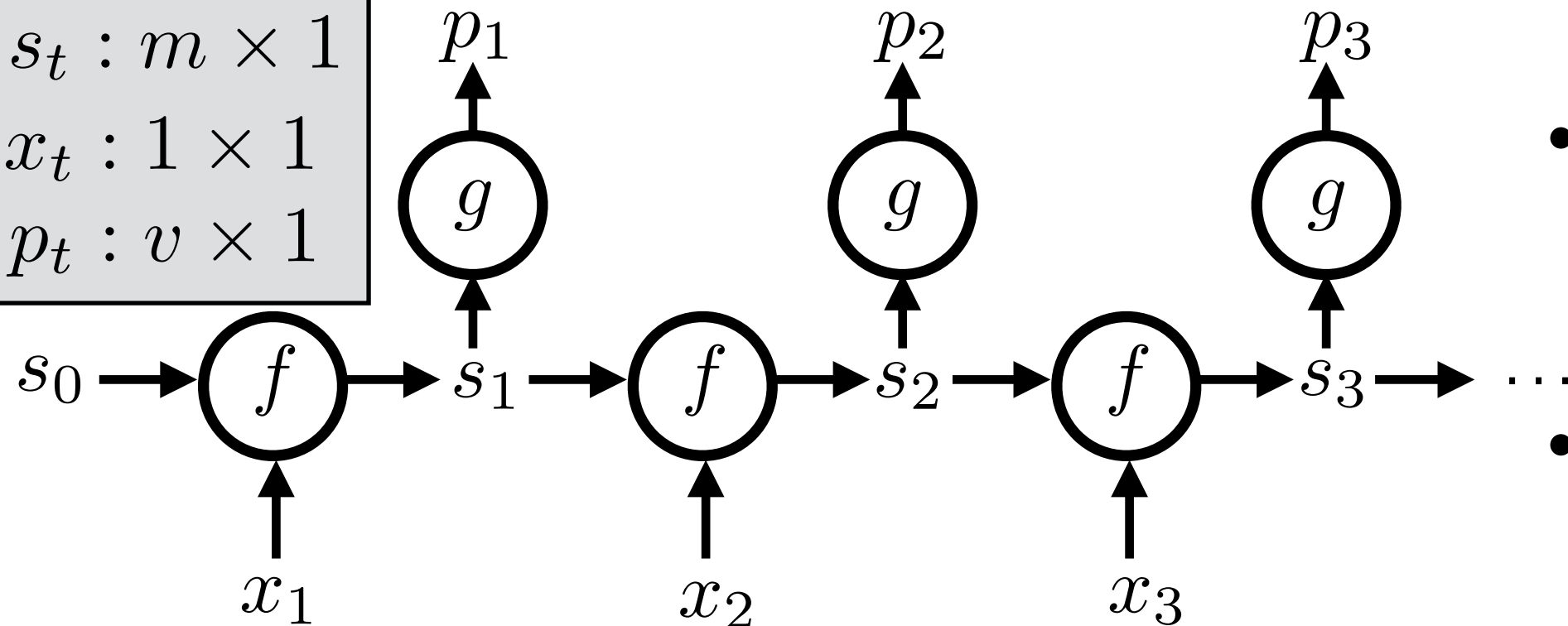


- Example: Alphabet $\{0, 1\}$;
state is last $m = 3$ characters

- m : number of characters in the context
- v : number of characters in the alphabet

Can express as a state machine

$$\begin{array}{l} s_t : m \times 1 \\ x_t : 1 \times 1 \\ p_t : v \times 1 \end{array}$$



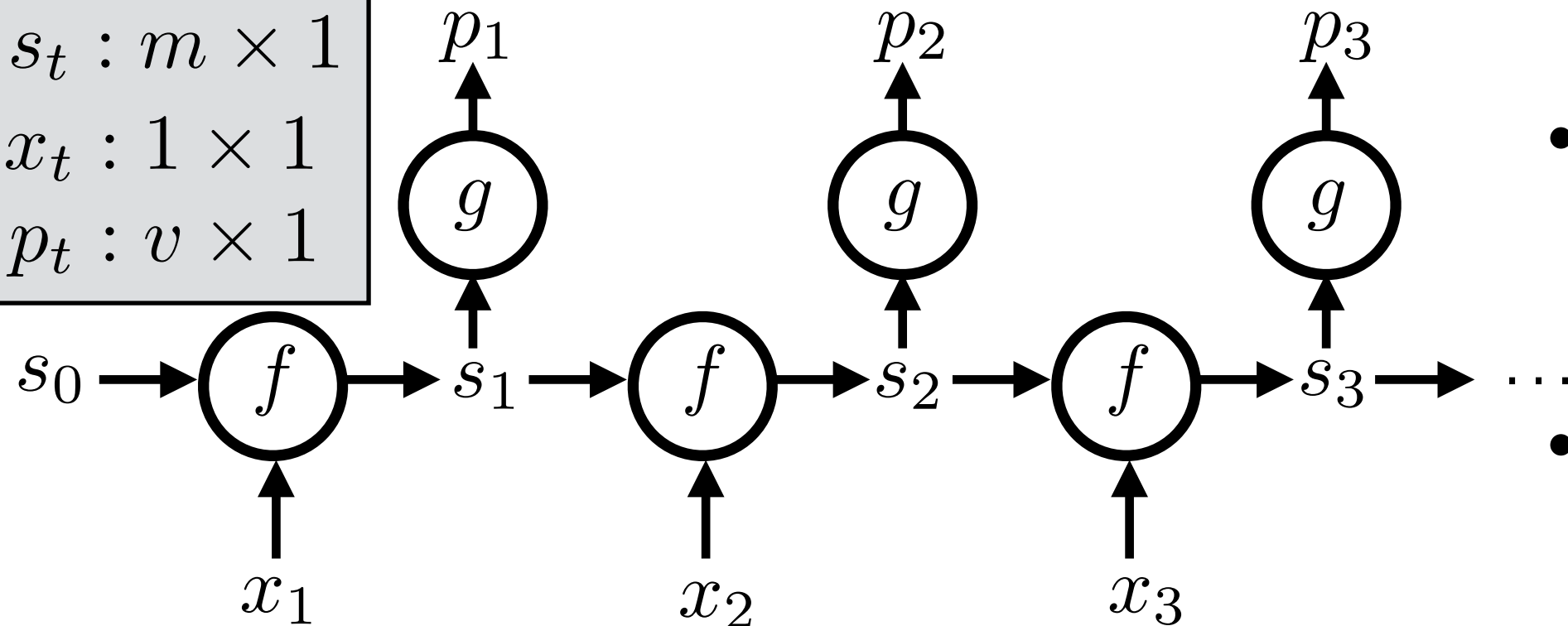
- m : number of characters in the context
- v : number of characters in the alphabet

- Example: Alphabet $\{0, 1\}$;
state is last $m = 3$ characters

$$s_t = f(s_{t-1}, x_t) =$$

Can express as a state machine

$$\begin{array}{l} s_t : m \times 1 \\ x_t : 1 \times 1 \\ p_t : v \times 1 \end{array}$$



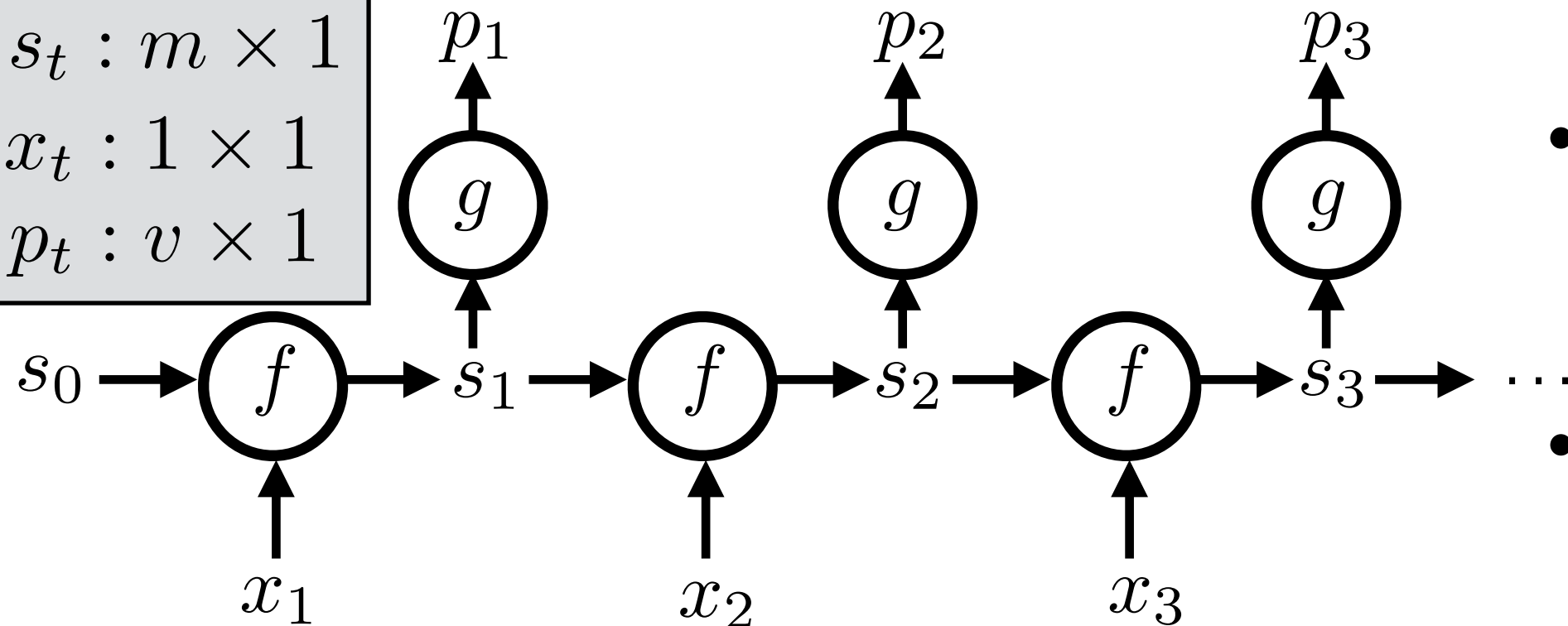
- Example: Alphabet $\{0, 1\}$;
state is last $m = 3$ characters

- m : number of characters in the context
- v : number of characters in the alphabet

$$s_t = f(s_{t-1}, x_t) =$$

Can express as a state machine

$$\begin{array}{l} s_t : m \times 1 \\ x_t : 1 \times 1 \\ p_t : v \times 1 \end{array}$$



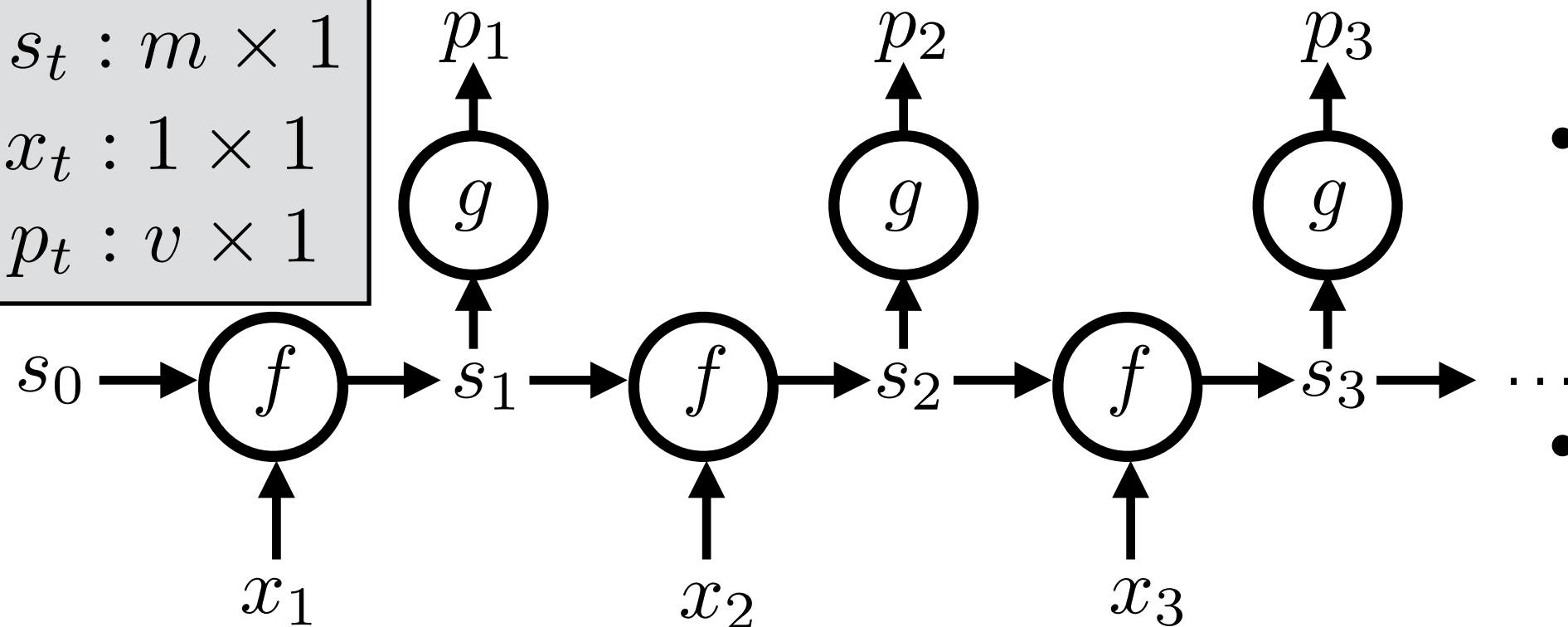
- m : number of characters in the context
- v : number of characters in the alphabet

- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = f(s_{t-1}, x_t) = \boxed{?} x_t + \boxed{?} s_{t-1}$$

Can express as a state machine

$$\begin{array}{l} s_t : m \times 1 \\ x_t : 1 \times 1 \\ p_t : v \times 1 \end{array}$$



- m : number of characters in the context
- v : number of characters in the alphabet

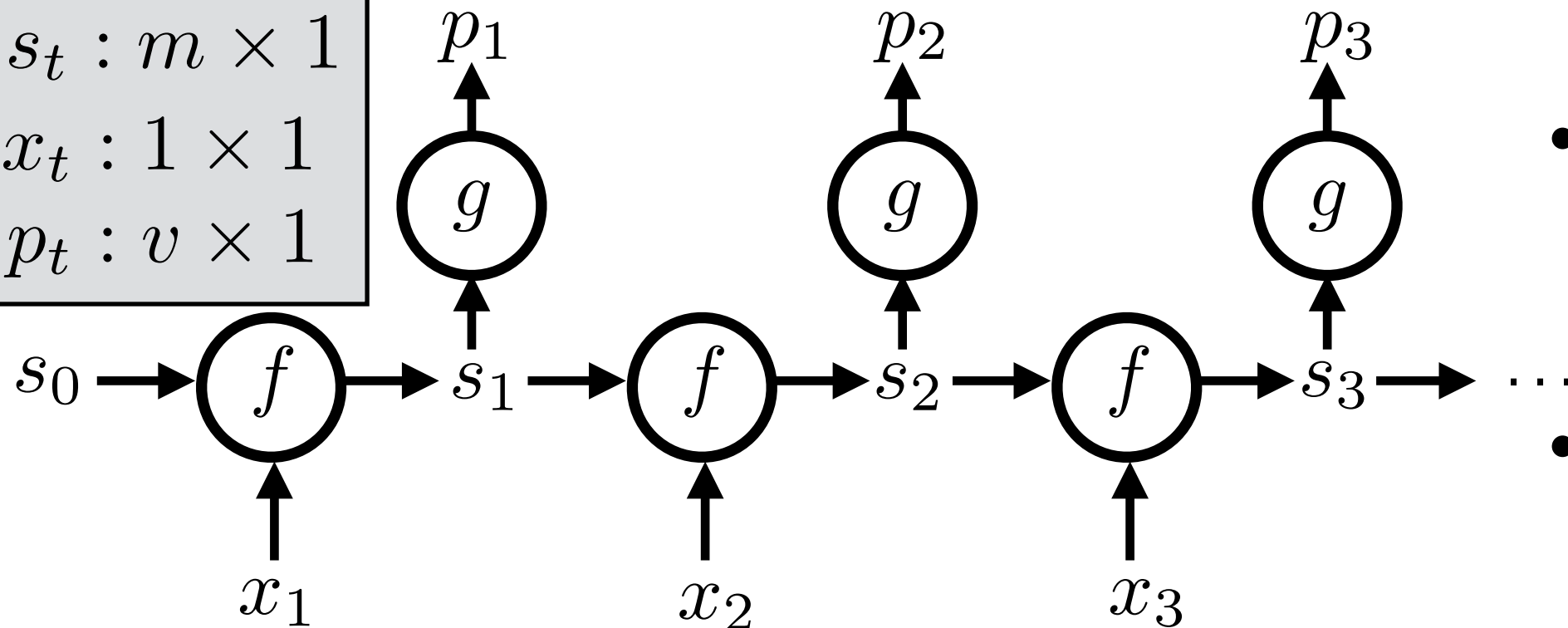
- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = f(s_{t-1}, x_t) = \begin{array}{|c|} \hline ? \\ \hline \end{array} x_t + \begin{array}{|c|} \hline ? \\ \hline \end{array} s_{t-1}$$

3×1

Can express as a state machine

$$\begin{array}{l} s_t : m \times 1 \\ x_t : 1 \times 1 \\ p_t : v \times 1 \end{array}$$



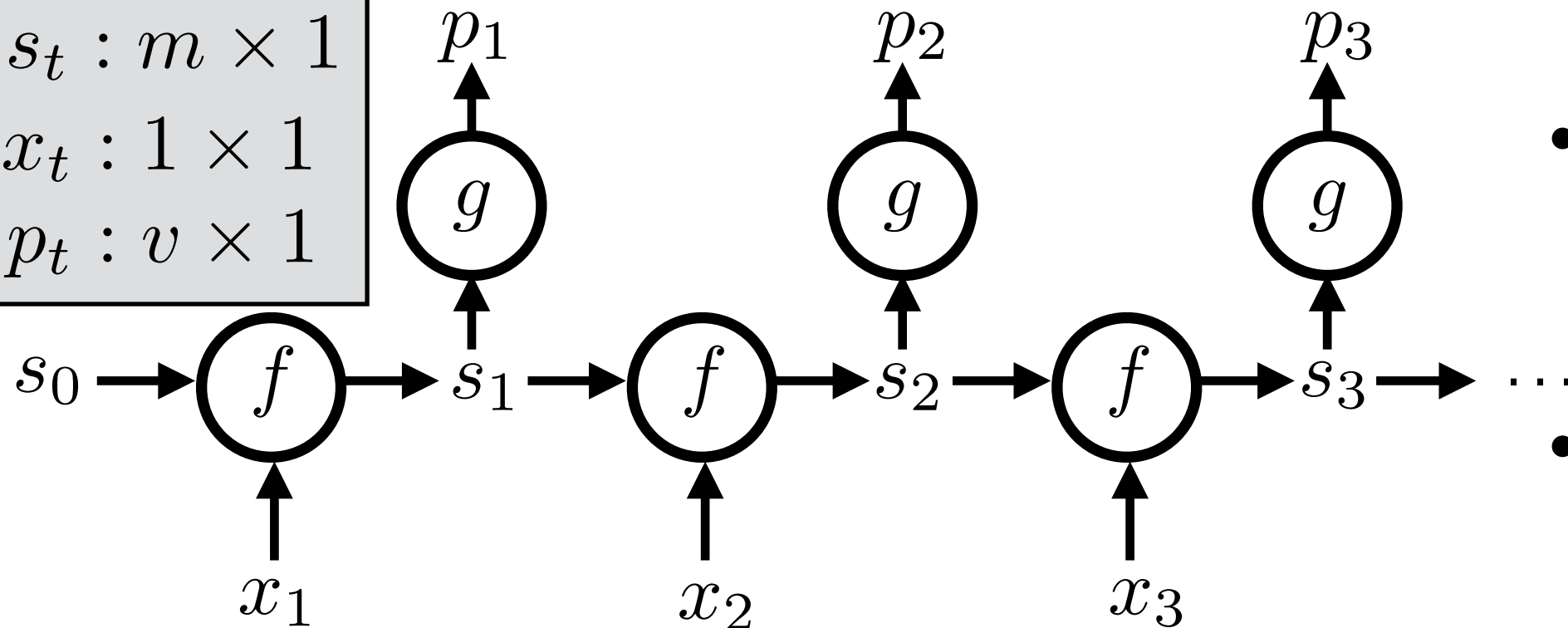
- m : number of characters in the context
- v : number of characters in the alphabet

- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = f(s_{t-1}, x_t) = \begin{array}{c} \text{3} \times \text{1} \\ \boxed{?} \end{array} x_t + \begin{array}{c} \boxed{?} \\ \text{1} \times \text{1} \end{array} s_{t-1}$$

Can express as a state machine

$$\begin{array}{l} s_t : m \times 1 \\ x_t : 1 \times 1 \\ p_t : v \times 1 \end{array}$$



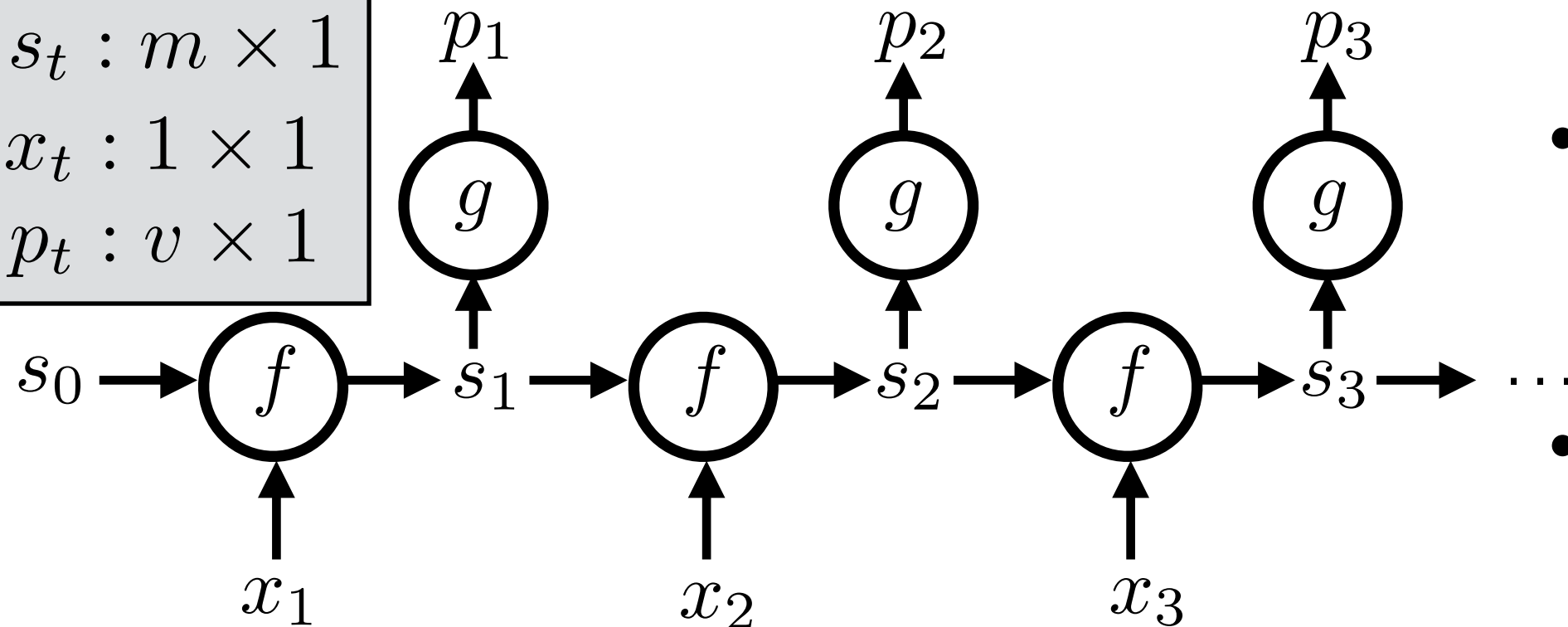
- m : number of characters in the context
- v : number of characters in the alphabet

- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = f(s_{t-1}, x_t) = \begin{array}{c} \text{?} \\ 3 \times 1 \end{array} x_t + \begin{array}{c} \text{?} \\ 3 \times 1, 1 \times 1 \end{array} s_{t-1}$$

Can express as a state machine

$$\begin{array}{l} s_t : m \times 1 \\ x_t : 1 \times 1 \\ p_t : v \times 1 \end{array}$$



- m : number of characters in the context
- v : number of characters in the alphabet

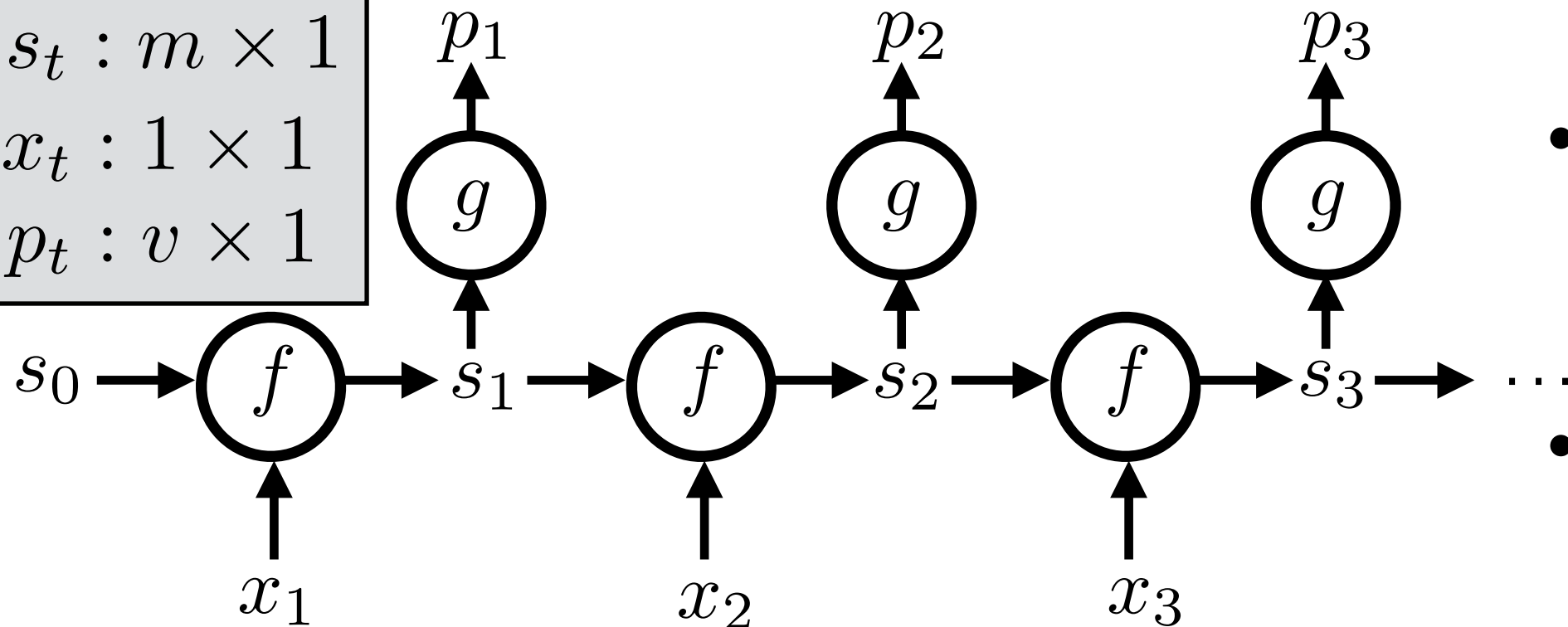
- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = f(s_{t-1}, x_t) = \begin{array}{c} \text{3} \times \text{1} \\ \text{?} \end{array} x_t + \begin{array}{c} \text{?} \end{array} s_{t-1}$$

$3 \times 1, 1 \times 1$

Can express as a state machine

$$\begin{array}{l} s_t : m \times 1 \\ x_t : 1 \times 1 \\ p_t : v \times 1 \end{array}$$



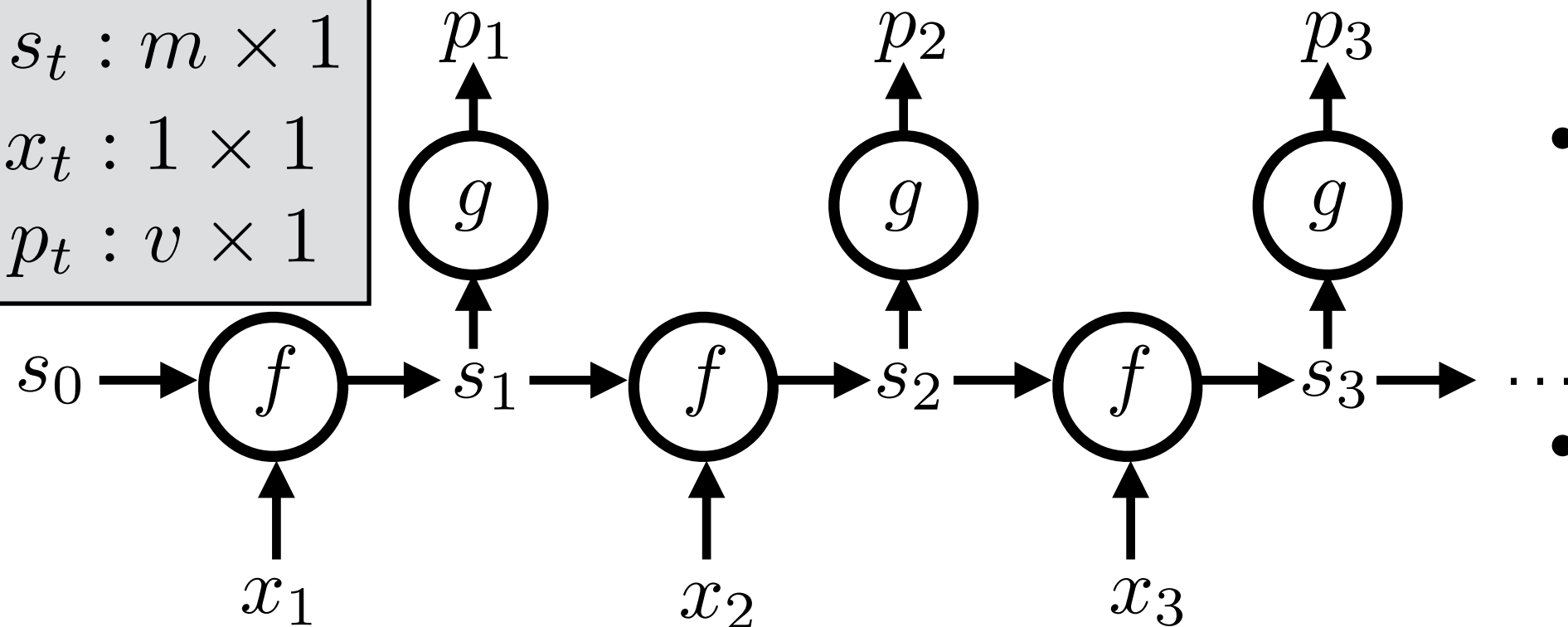
- m : number of characters in the context
- v : number of characters in the alphabet

- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} x_t + \begin{array}{c} \text{?} \end{array} s_{t-1}$$

Can express as a state machine

$$\begin{array}{l} s_t : m \times 1 \\ x_t : 1 \times 1 \\ p_t : v \times 1 \end{array}$$



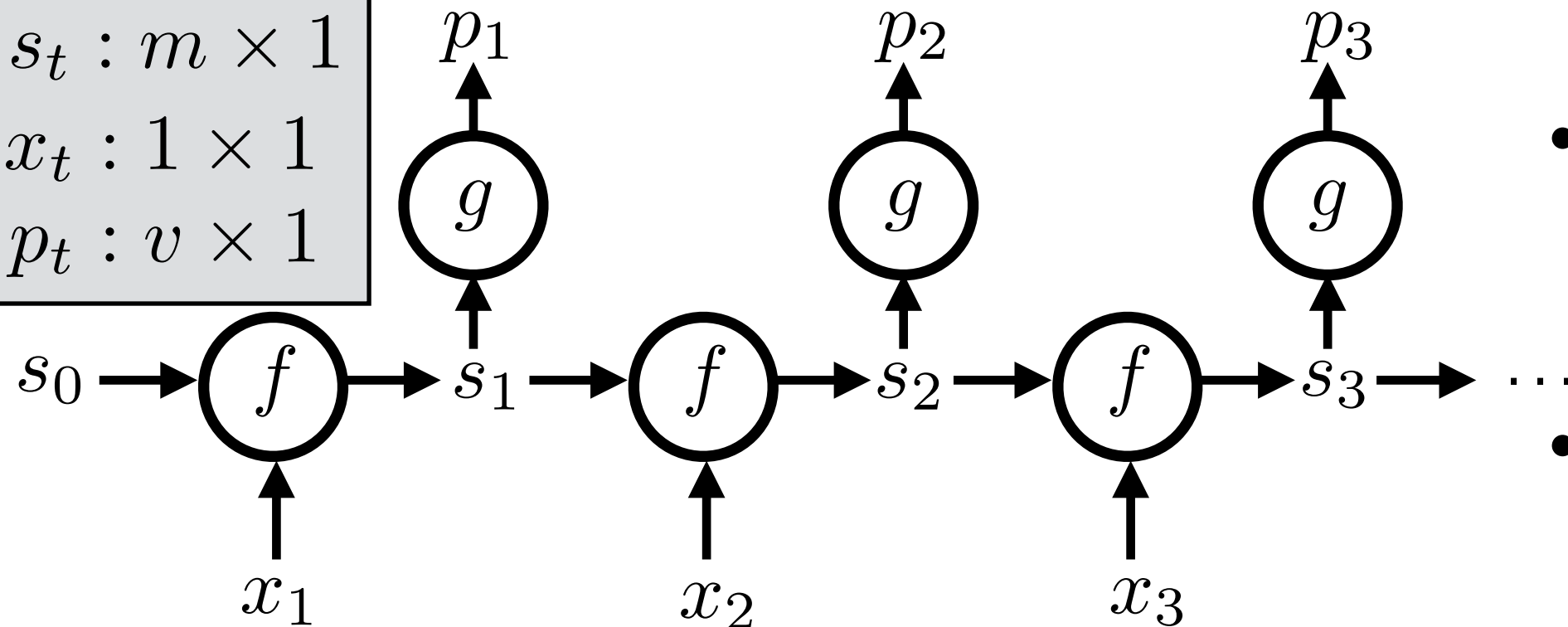
- m : number of characters in the context
- v : number of characters in the alphabet

- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} x_t + \begin{bmatrix} ? \\ ? \\ ? \end{bmatrix} s_{t-1}$$

Can express as a state machine

$$\begin{array}{l} s_t : m \times 1 \\ x_t : 1 \times 1 \\ p_t : v \times 1 \end{array}$$



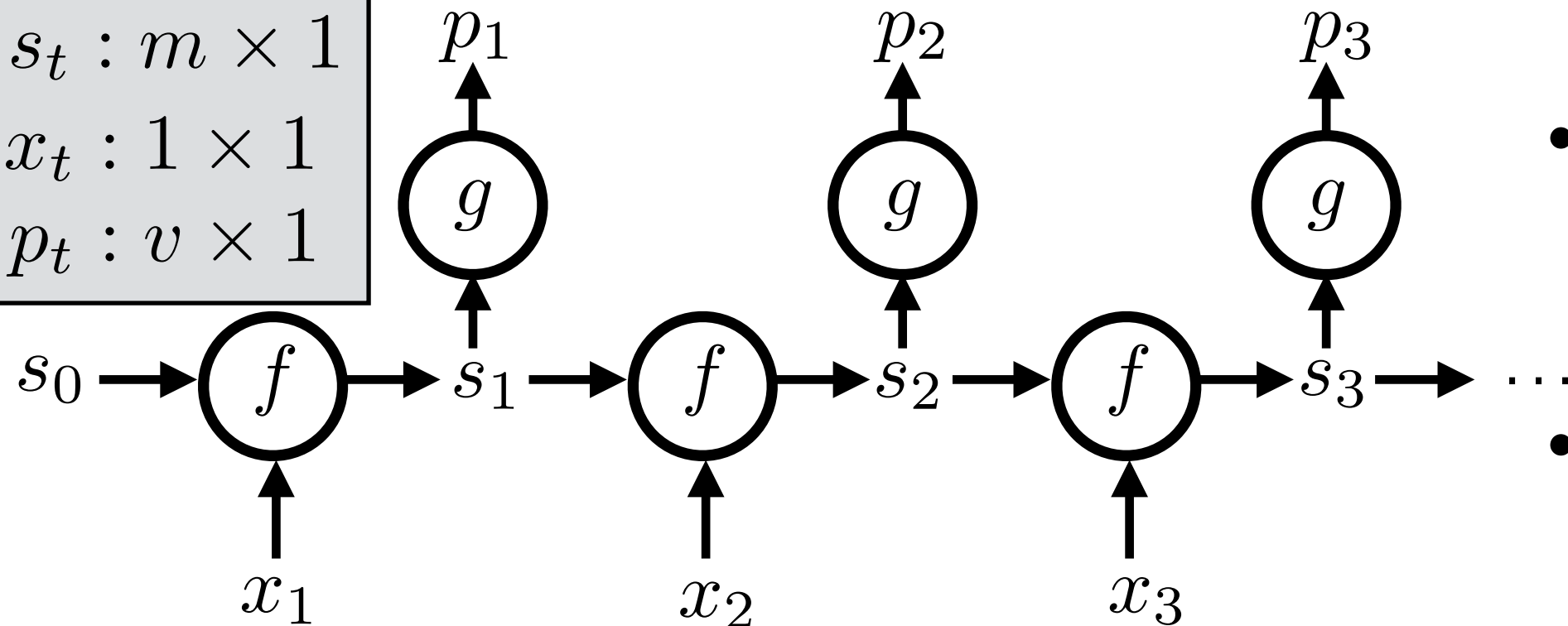
- m : number of characters in the context
- v : number of characters in the alphabet

- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} x_t + \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} s_{t-1}$$

Can express as a state machine

$$\begin{array}{l} s_t : m \times 1 \\ x_t : 1 \times 1 \\ p_t : v \times 1 \end{array}$$



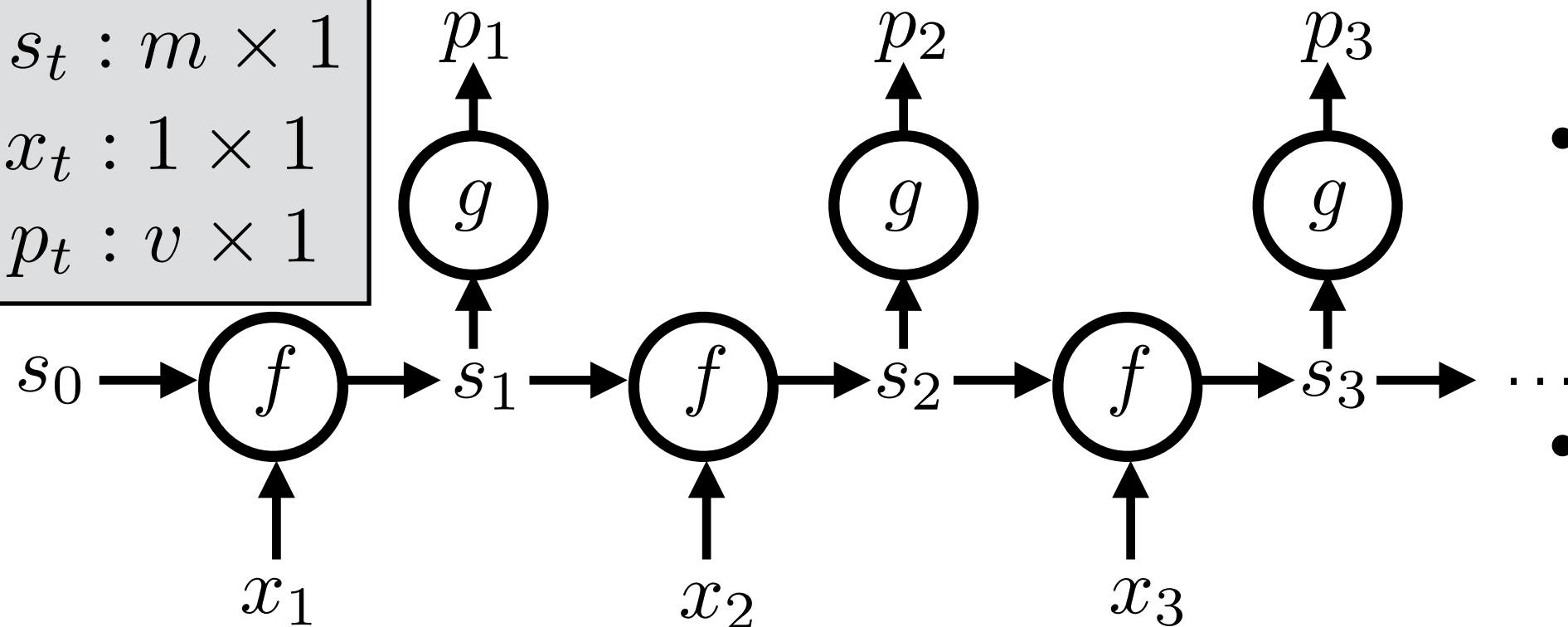
- m : number of characters in the context
- v : number of characters in the alphabet

- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} x_t + \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} s_{t-1}$$

Can express as a state machine

$$\begin{array}{l} s_t : m \times 1 \\ x_t : 1 \times 1 \\ p_t : v \times 1 \end{array}$$



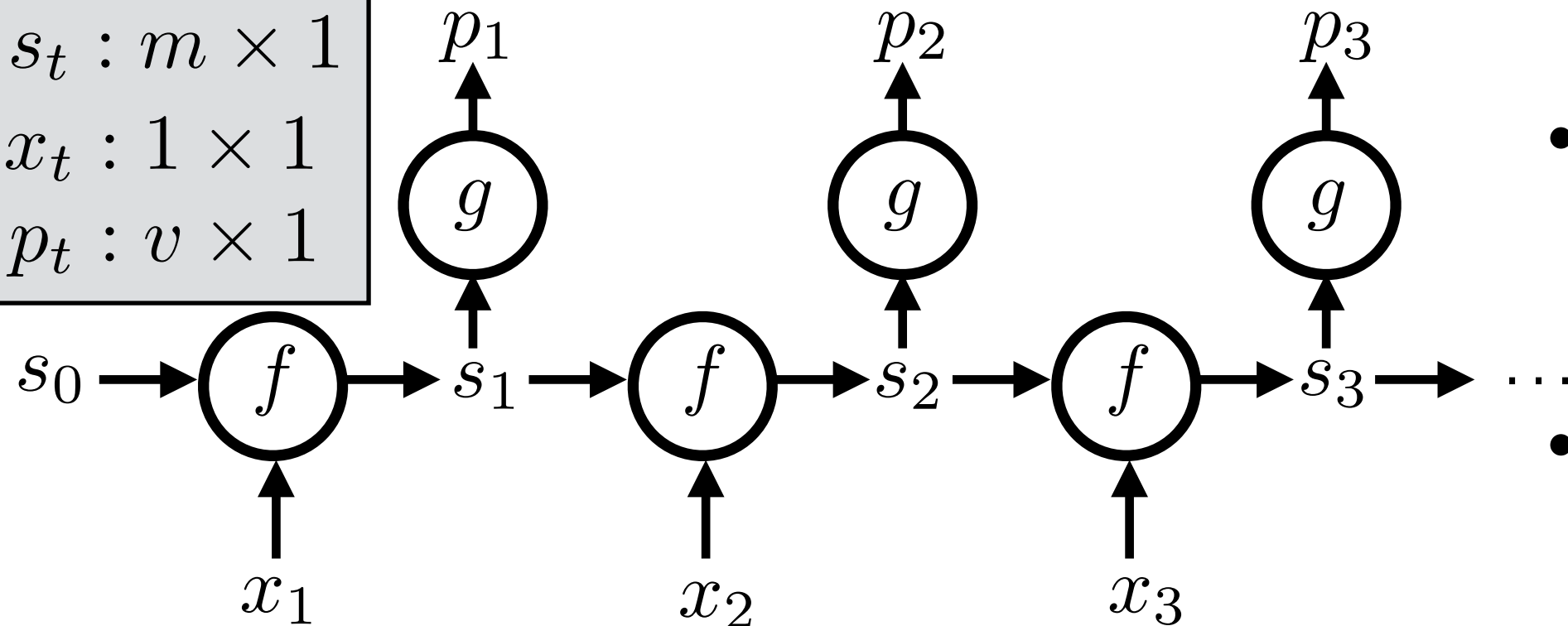
- m : number of characters in the context
- v : number of characters in the alphabet

- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} x_t + \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} s_{t-1}$$

Can express as a state machine

$$\begin{array}{l} s_t : m \times 1 \\ x_t : 1 \times 1 \\ p_t : v \times 1 \end{array}$$



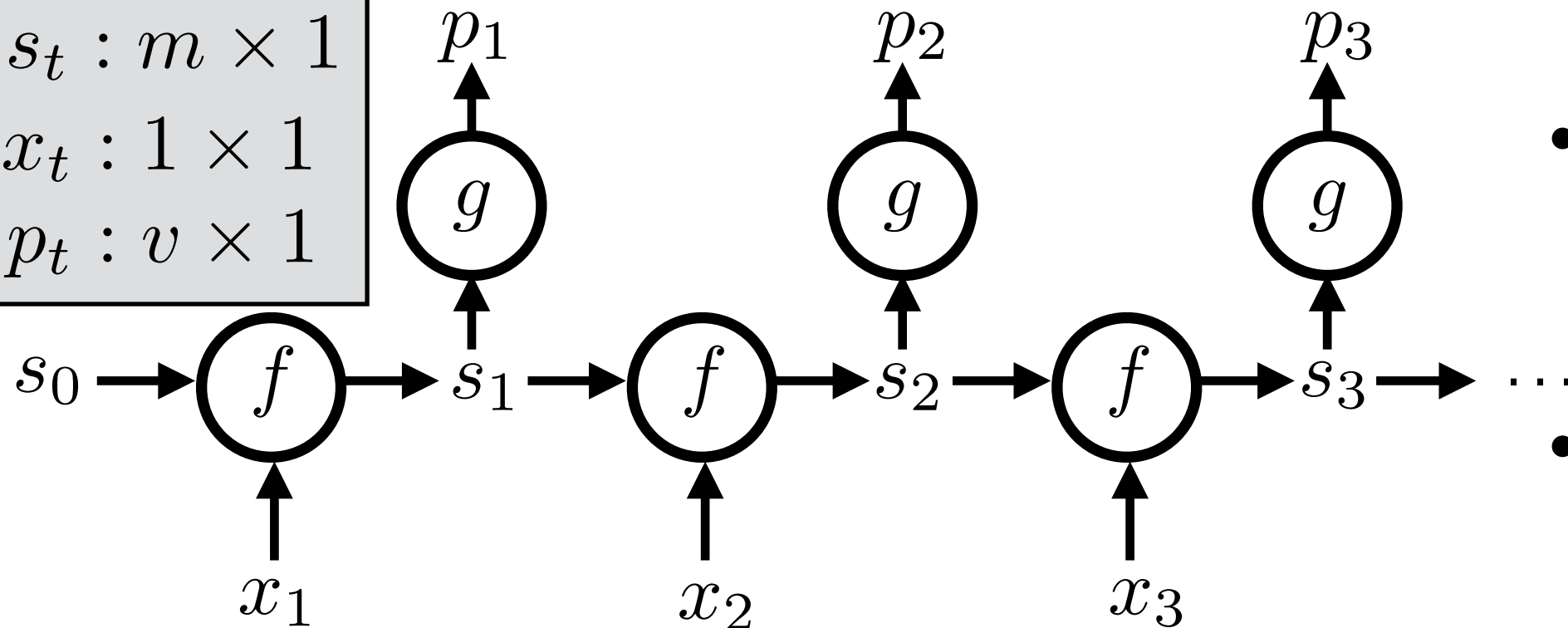
- m : number of characters in the context
- v : number of characters in the alphabet

- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} x_t + \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} s_{t-1}$$

Can express as a state machine

$$\begin{array}{l} s_t : m \times 1 \\ x_t : 1 \times 1 \\ p_t : v \times 1 \end{array}$$



- m : number of characters in the context
- v : number of characters in the alphabet

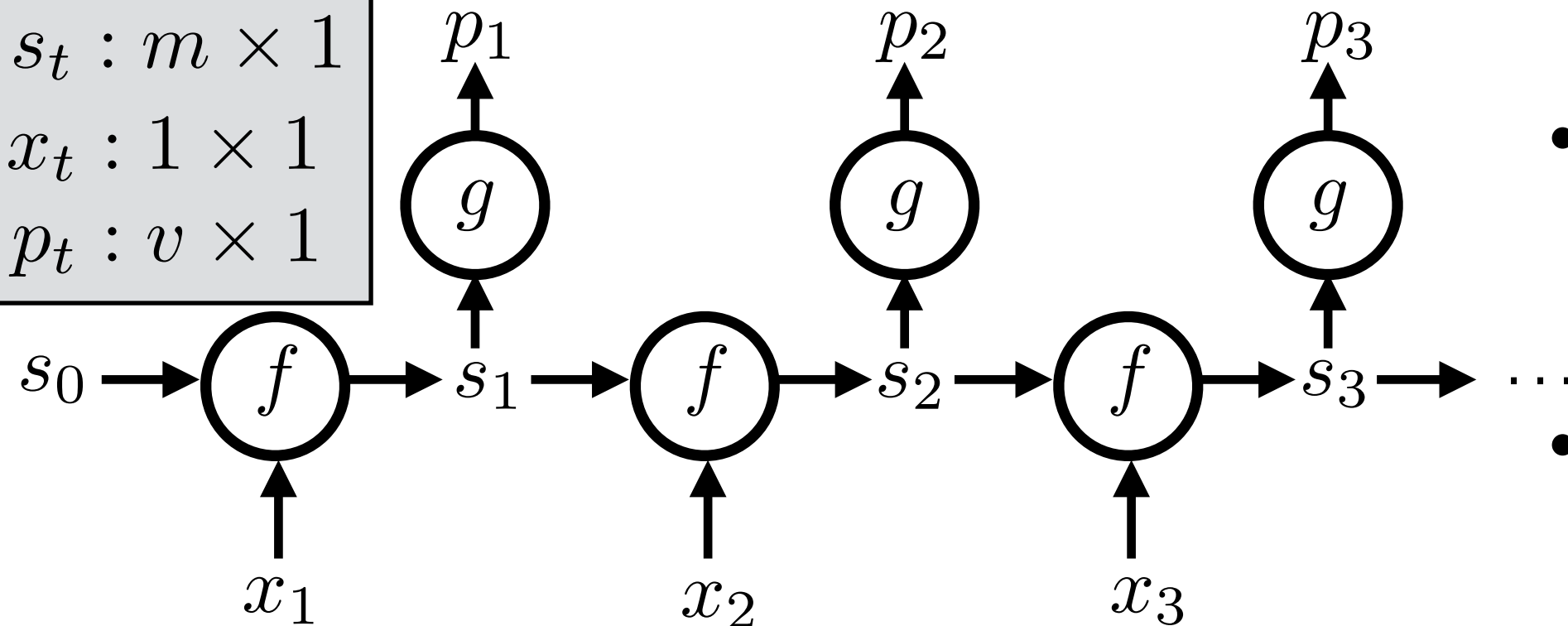
- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} x_t + \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} s_{t-1}$$

$$p_t = g(s_t)$$

Can express as a state machine

$$\begin{array}{l} s_t : m \times 1 \\ x_t : 1 \times 1 \\ p_t : v \times 1 \end{array}$$



- m : number of characters in the context
- v : number of characters in the alphabet

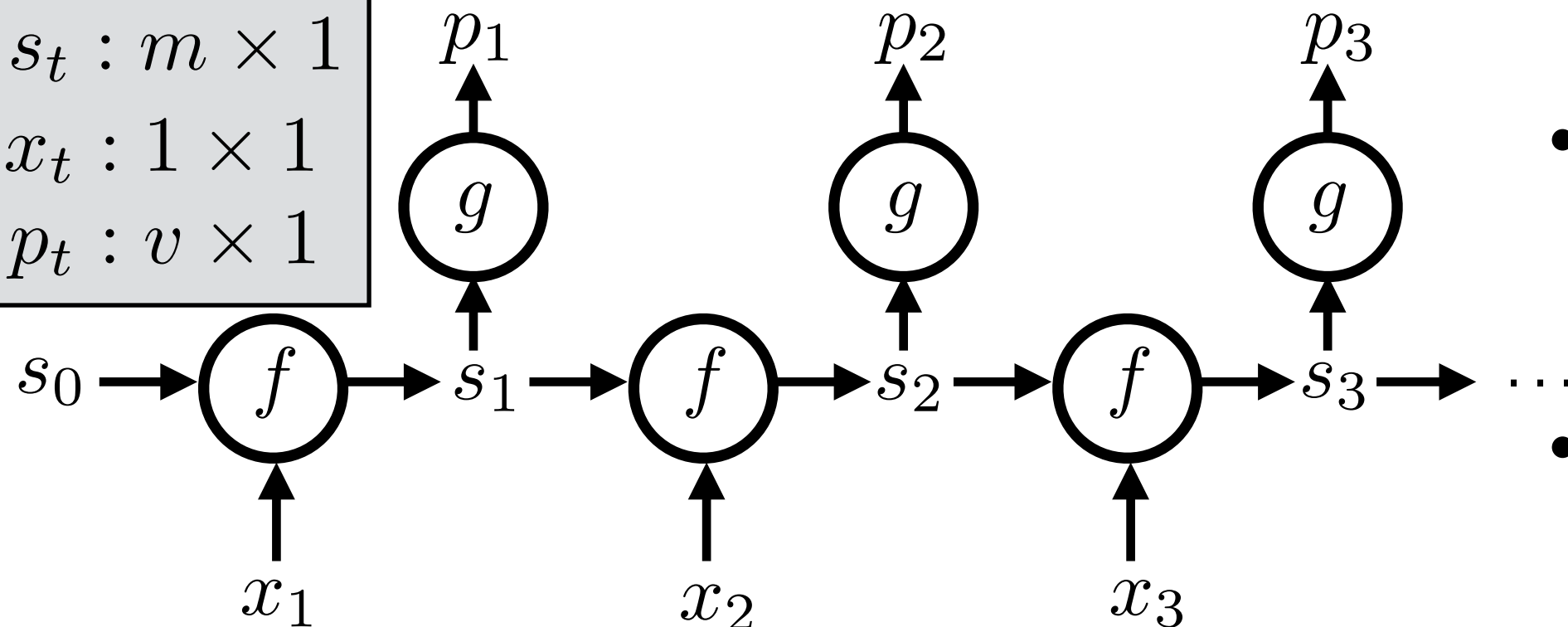
- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} x_t + \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} s_{t-1}$$

$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(W^o s_t + W_0^o) \end{aligned}$$

Can express as a state machine

$$\begin{array}{l} s_t : m \times 1 \\ x_t : 1 \times 1 \\ p_t : v \times 1 \end{array}$$



- m : number of characters in the context
- v : number of characters in the alphabet

- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

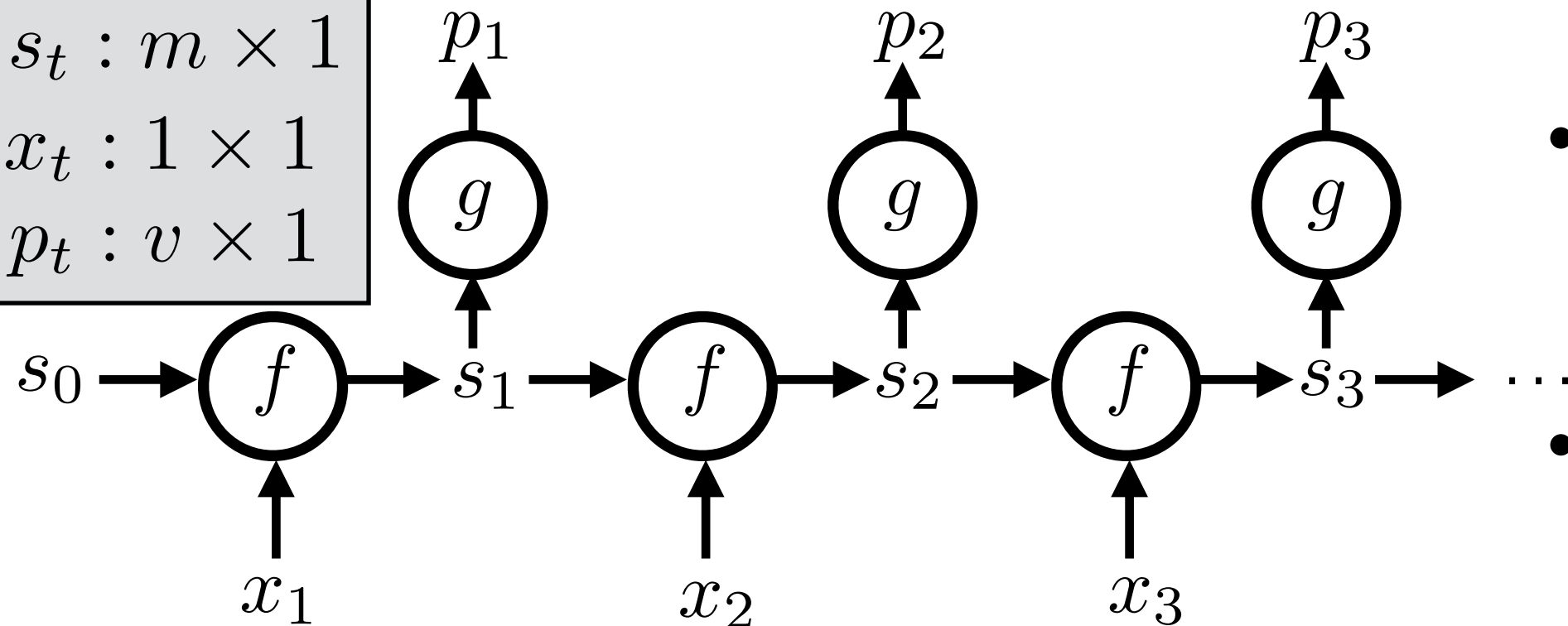
$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} x_t + \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} s_{t-1}$$

$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(W^o s_t + W_0^o) \end{aligned}$$

2-class
logistic
regression

Can express as a state machine

$$\begin{array}{l} s_t : m \times 1 \\ x_t : 1 \times 1 \\ p_t : v \times 1 \end{array}$$



- m : number of characters in the context
- v : number of characters in the alphabet

- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} x_t + \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} s_{t-1}$$

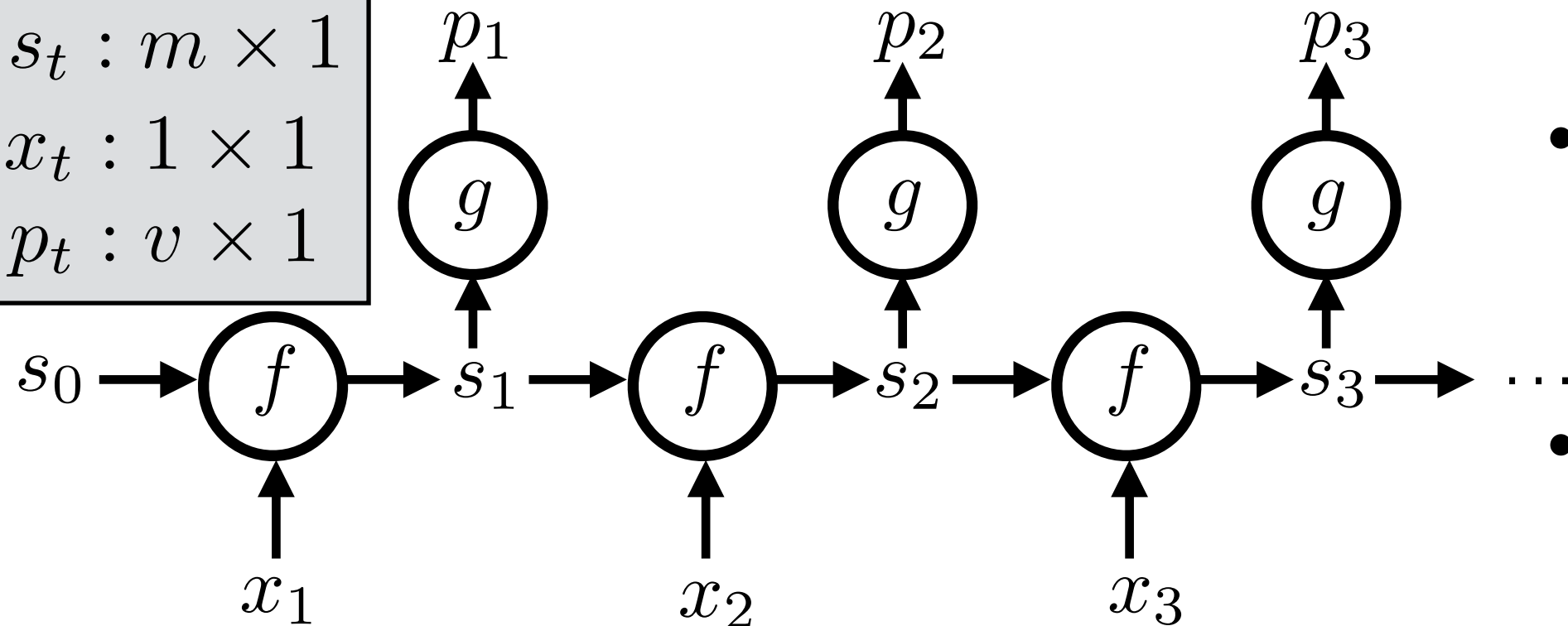
$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(W^o s_t + W_0^o) \end{aligned}$$

1×3

2-class
logistic
regression

Can express as a state machine

$$\begin{array}{l} s_t : m \times 1 \\ x_t : 1 \times 1 \\ p_t : v \times 1 \end{array}$$



- m : number of characters in the context
- v : number of characters in the alphabet

- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} x_t + \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} s_{t-1}$$

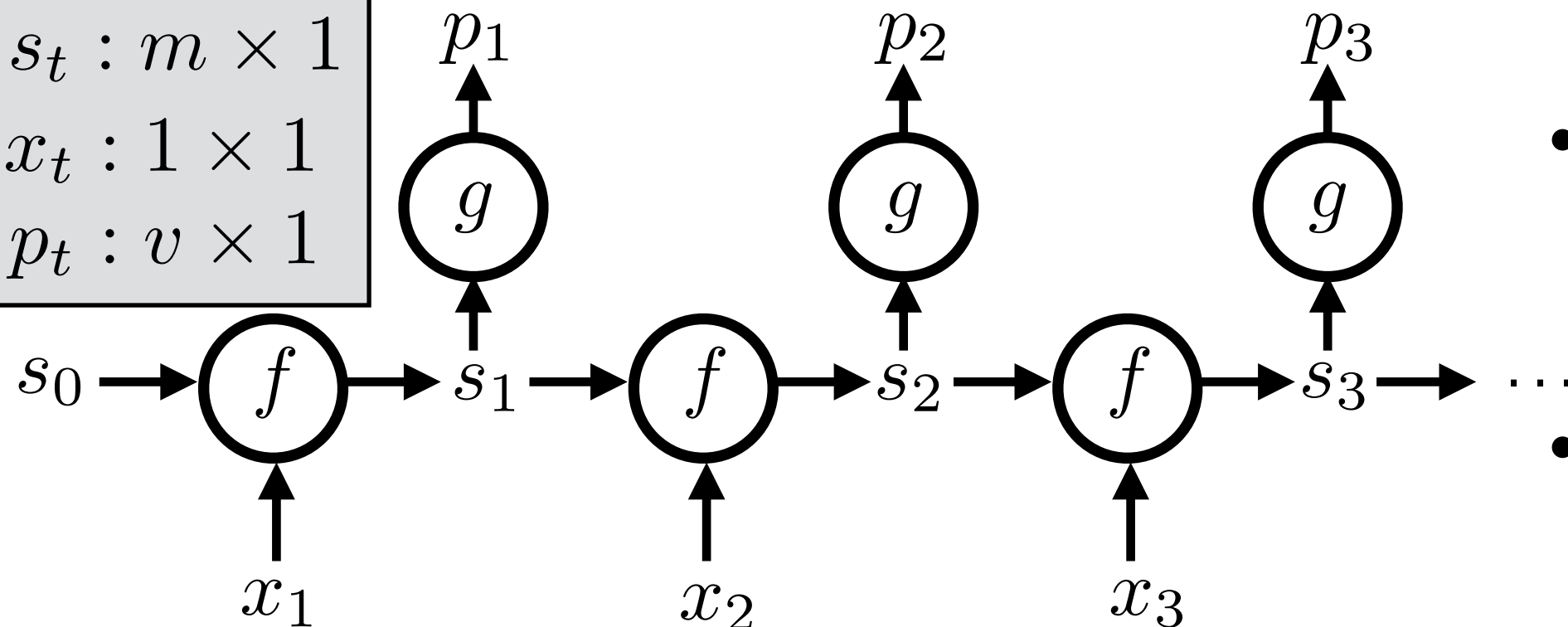
$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(W^o s_t + W_0^o) \end{aligned}$$

1×3 1×1

2-class
logistic
regression

Can express as a state machine

$$\begin{array}{l} s_t : m \times 1 \\ x_t : 1 \times 1 \\ p_t : v \times 1 \end{array}$$



- m : number of characters in the context
- v : number of characters in the alphabet

- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

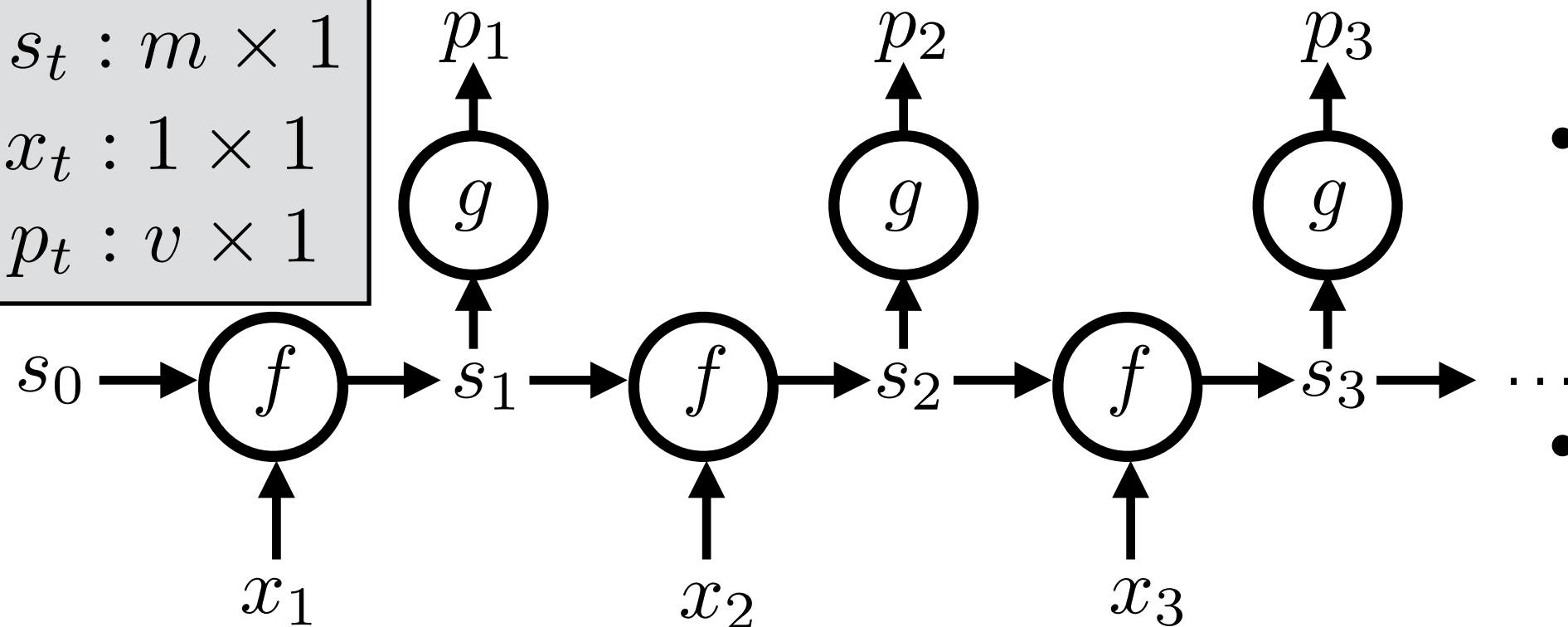
$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} x_t + \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} s_{t-1}$$

$$\begin{aligned} p_t &= g(s_t) \\ &= \underset{1 \times 3}{f_2}(\underset{1 \times 1}{W^o} s_t + W_0^o) \end{aligned}$$

2-class
logistic
regression

Can express as a state machine

$$\begin{array}{l} s_t : m \times 1 \\ x_t : 1 \times 1 \\ p_t : v \times 1 \end{array}$$



- m : number of characters in the context
- v : number of characters in the alphabet

- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} x_t + \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} s_{t-1}$$

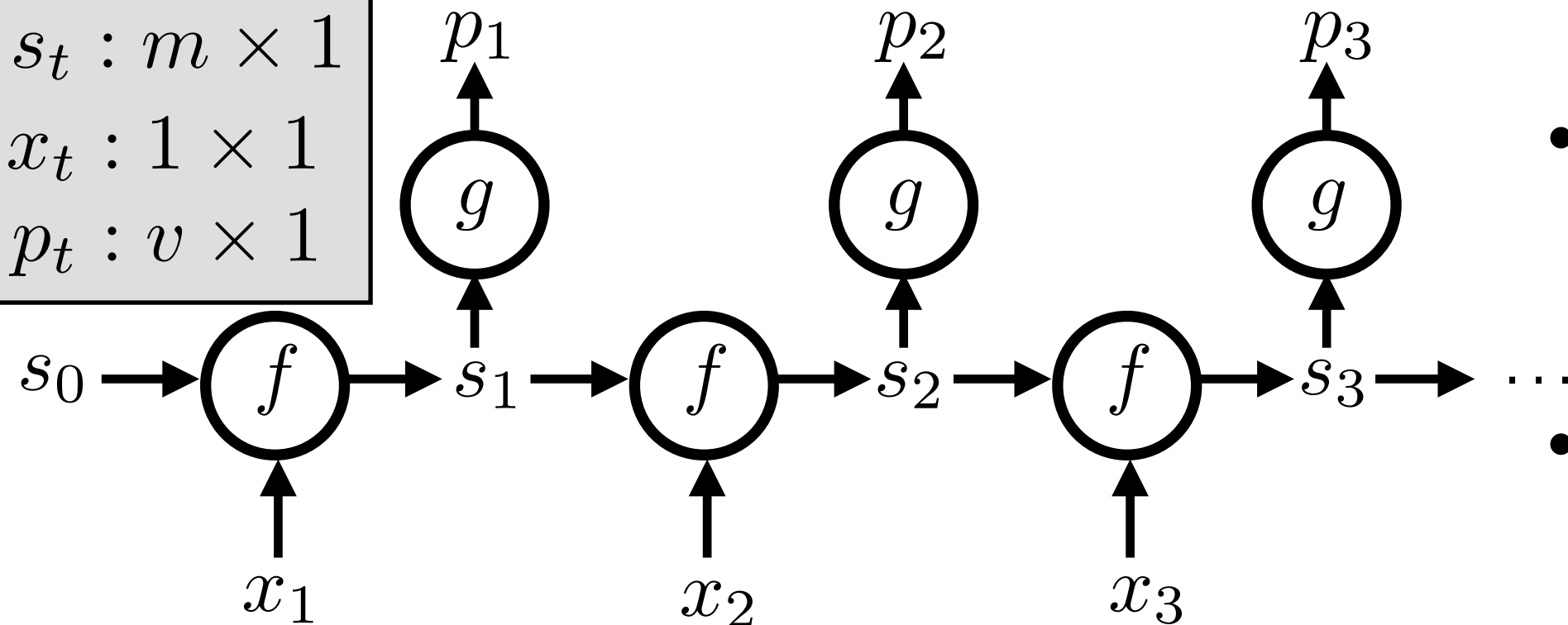
$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(W^o s_t + W_0^o) \end{aligned}$$

$1 \times 3 \quad 1 \times 1$

v -class
logistic
regression

Can express as a state machine

$$\begin{matrix} s_t : m \times 1 \\ x_t : 1 \times 1 \\ p_t : v \times 1 \end{matrix}$$



- m : number of characters in the context
- v : number of characters in the alphabet

- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

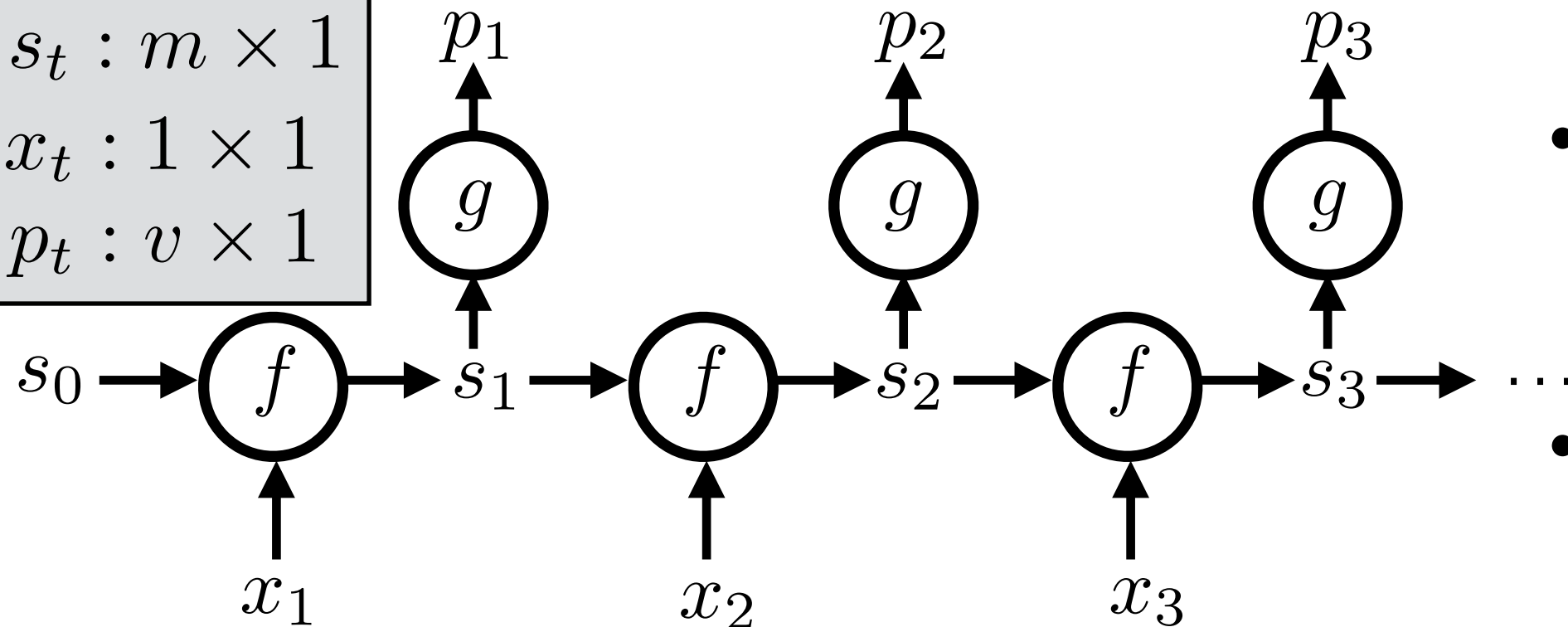
$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} x_t + \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} s_{t-1}$$

$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(\underbrace{W^o}_{1 \times 3} s_t + \underbrace{W_0^o}_{1 \times 1}) \end{aligned}$$

v -class
logistic
regression

Can express as a state machine

$$\begin{array}{l} s_t : m \times 1 \\ x_t : 1 \times 1 \\ p_t : v \times 1 \end{array}$$



- m : number of characters in the context
- v : number of characters in the alphabet

- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

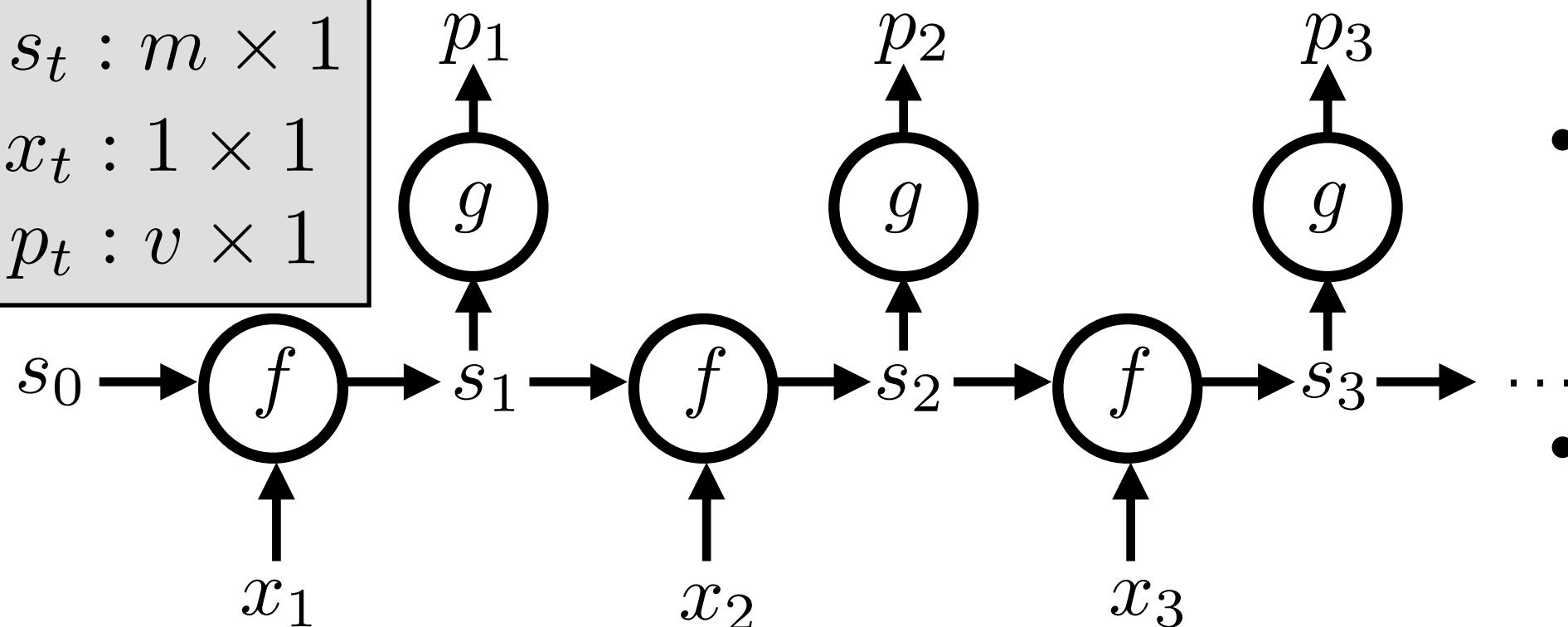
$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} x_t + \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} s_{t-1}$$

$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(\underbrace{W^o}_{v \times 3} s_t + \underbrace{W_0^o}_{v \times 1}) \end{aligned}$$

v -class
logistic
regression

Can express as a state machine

$$\begin{array}{l} s_t : m \times 1 \\ x_t : 1 \times 1 \\ p_t : v \times 1 \end{array}$$



- m : number of characters in the context
- v : number of characters in the alphabet

- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

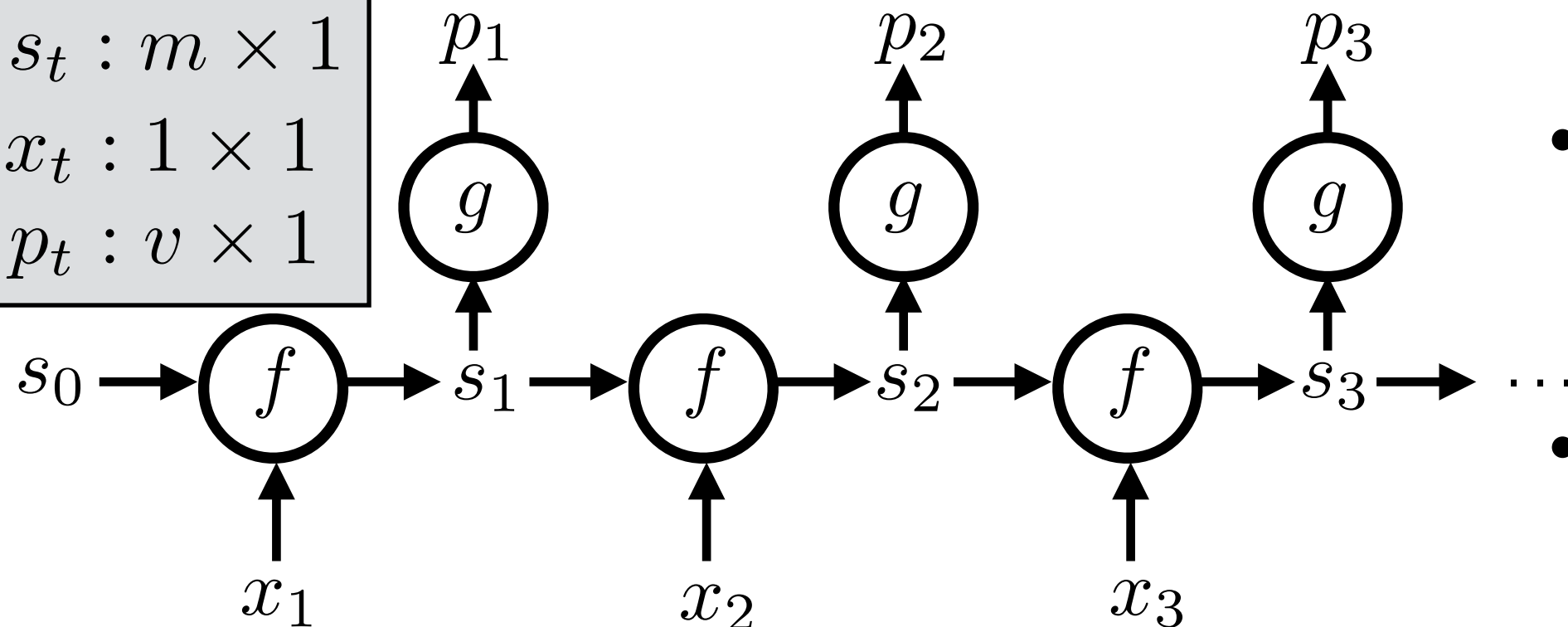
$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} x_t + \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} s_{t-1}$$

$$\begin{aligned} p_t &= g(s_t) \\ &= \underset{v \times 3}{f_2} (W^o s_t + \underset{v \times 1}{W_0^o}) \end{aligned}$$

v -class
logistic
regression

Can express as a state machine

$$\begin{array}{l} s_t : m \times 1 \\ x_t : 1 \times 1 \\ p_t : v \times 1 \end{array}$$



- m : number of characters in the context
- v : number of characters in the alphabet

- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} x_t + \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} s_{t-1}$$

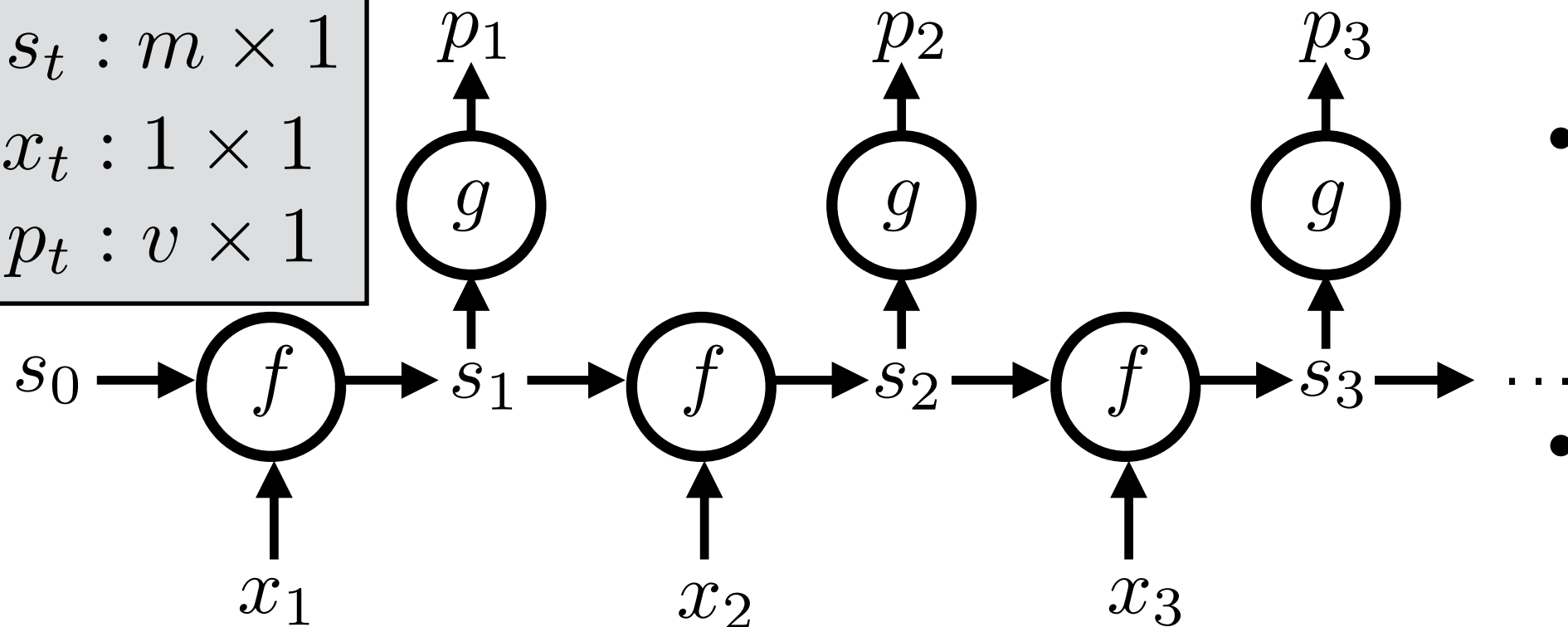
$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(W^o s_t + W_0^o) \end{aligned}$$

$v \times 3 \quad v \times 1$

v -class
logistic
regression

Can express as a state machine

$$\begin{array}{l} s_t : m \times 1 \\ x_t : 1 \times 1 \\ p_t : v \times 1 \end{array}$$



- m : number of characters in the context
- v : number of characters in the alphabet

- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} x_t + \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} s_{t-1}$$

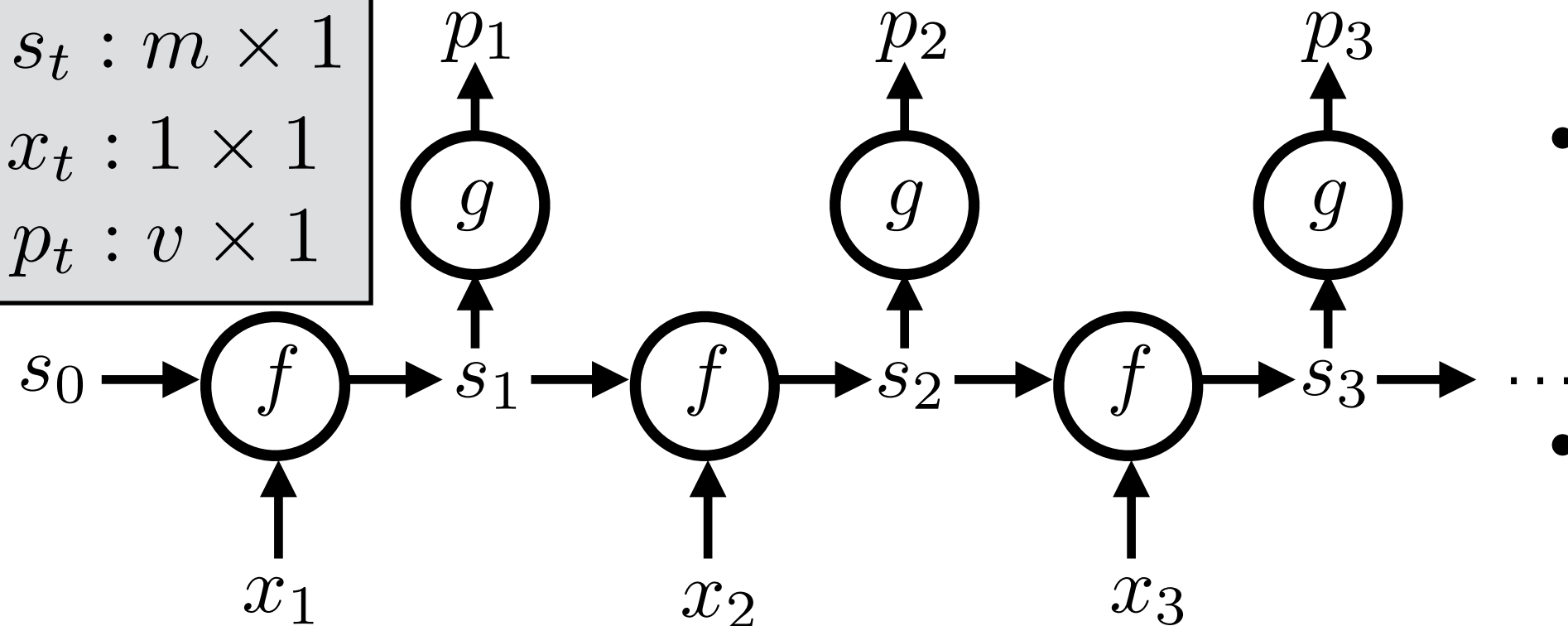
$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(W^o s_t + W_0^o) \end{aligned}$$

$v \times 3$ $v \times 1$

v -class
logistic
regression

Can express as a state machine

$$\begin{array}{l} s_t : m \times 1 \\ x_t : 1 \times 1 \\ p_t : v \times 1 \end{array}$$



- m : number of characters in the context
- v : number of characters in the alphabet

- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} x_t + \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} s_{t-1}$$

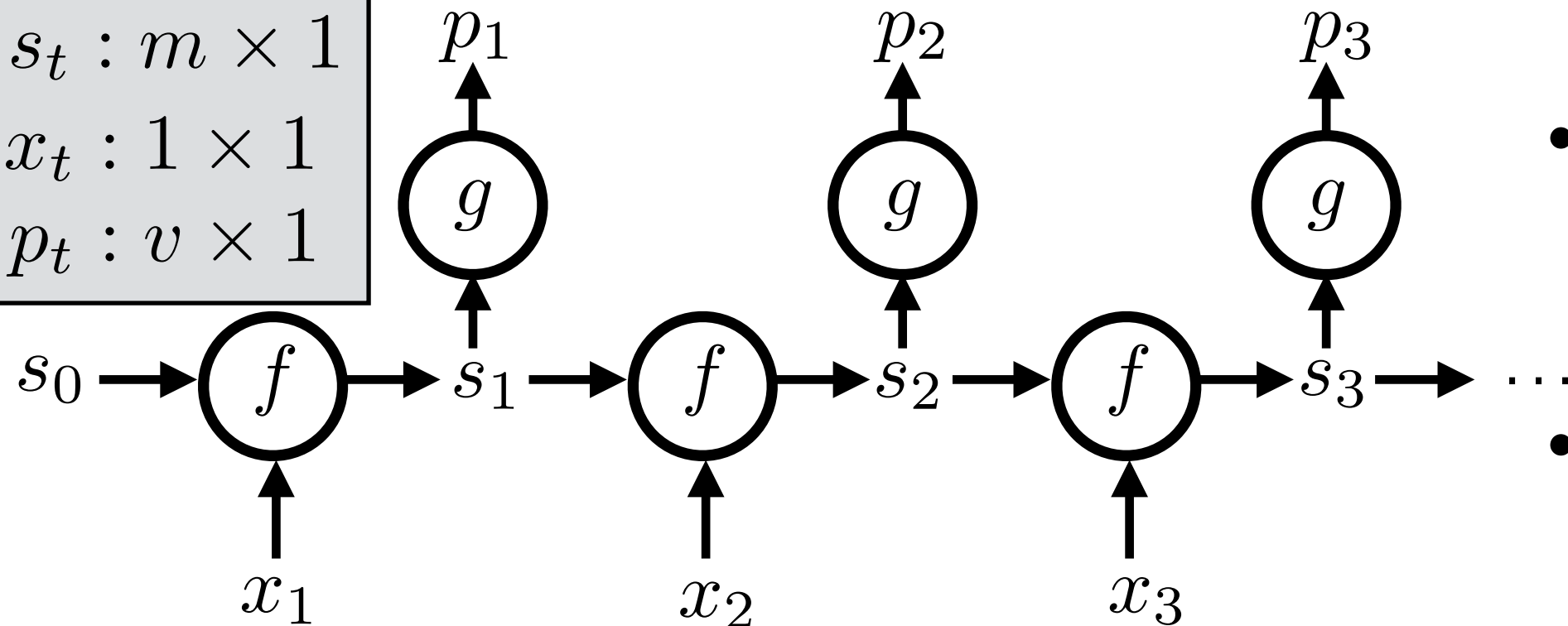
$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(W^o s_t + W_0^o) \end{aligned}$$

$v \times m$ $v \times 1$

v -class
logistic
regression

Can express as a state machine

$$\begin{array}{l} s_t : m \times 1 \\ x_t : 1 \times 1 \\ p_t : v \times 1 \end{array}$$



- m : number of characters in the context
- v : number of characters in the alphabet

- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} x_t + \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} s_{t-1}$$

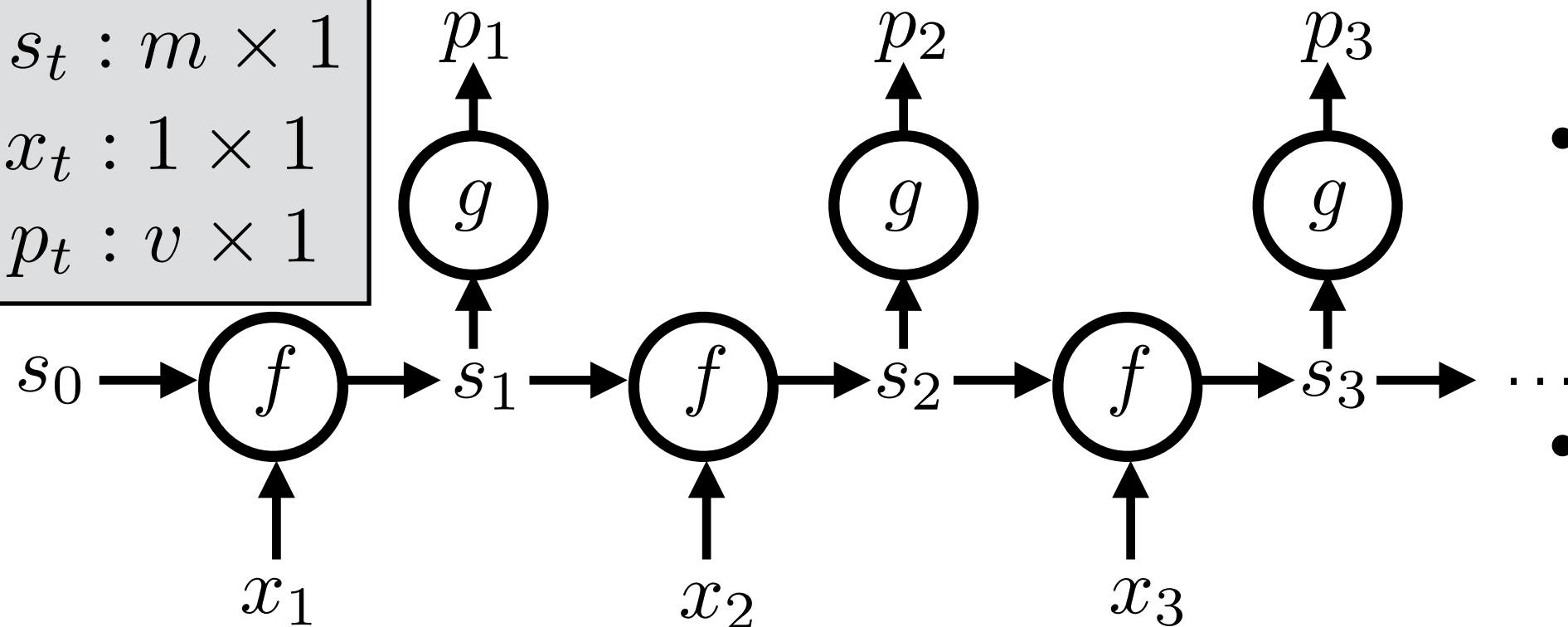
$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(W^o s_t + W_0^o) \end{aligned}$$

$v \times m \quad v \times 1$

v -class
logistic
regression

Can express as a state machine

$$\begin{array}{l} s_t : m \times 1 \\ x_t : 1 \times 1 \\ p_t : v \times 1 \end{array}$$



- m : number of characters in the context
- v : number of characters in the alphabet

- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

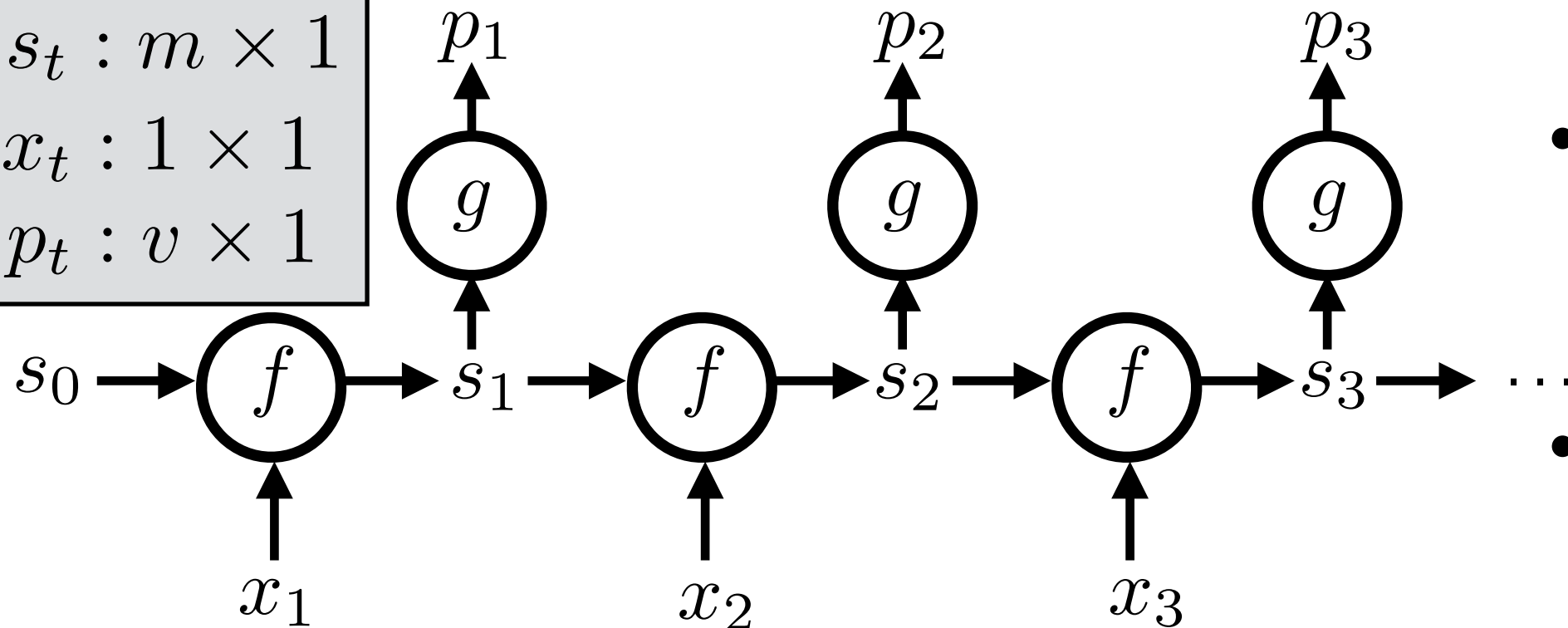
$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} x_t + \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} s_{t-1}$$

$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(\underbrace{W^o}_{v \times m} s_t + \underbrace{W_0^o}_{v \times 1}) \end{aligned}$$

v -class
logistic
regression

Can express as a state machine

$$\begin{matrix} s_t : m \times 1 \\ x_t : 1 \times 1 \\ p_t : v \times 1 \end{matrix}$$



- m : number of characters in the context
- v : number of characters in the alphabet

- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

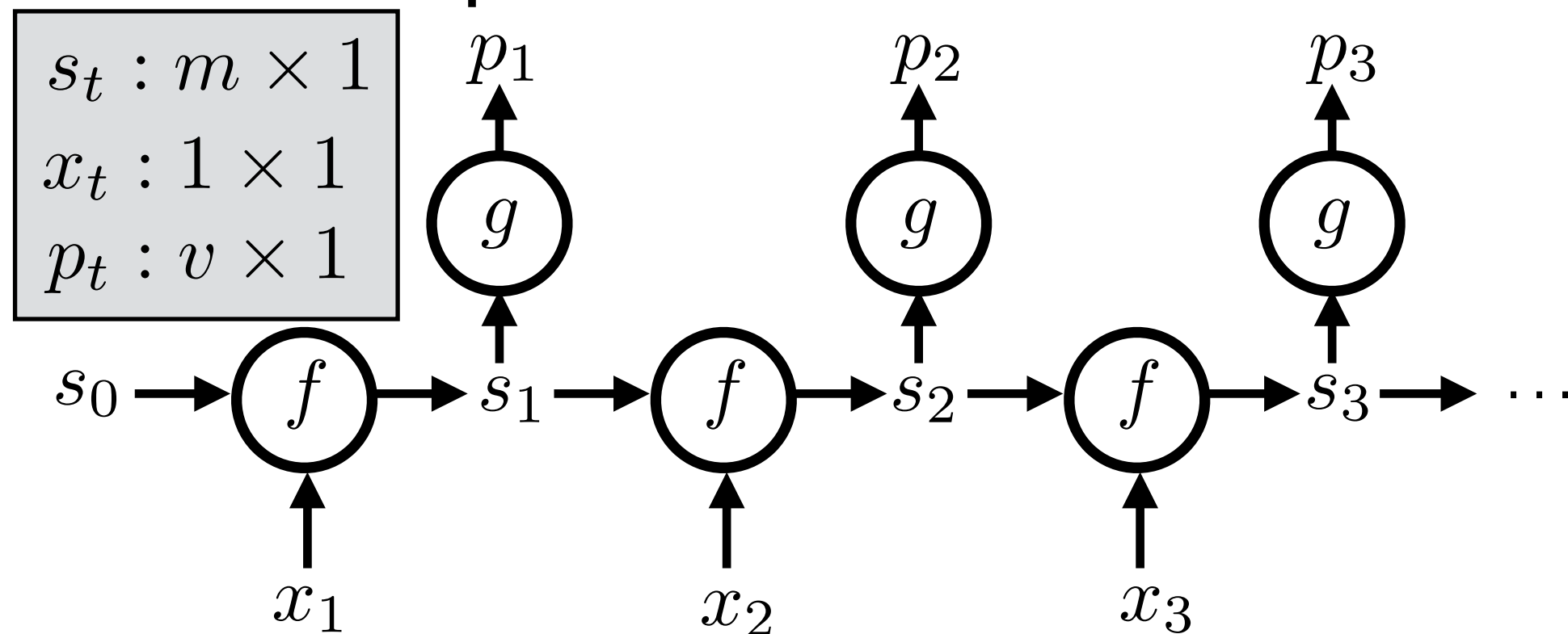
$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} x_t + \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} s_{t-1}$$

$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(\underbrace{W^o}_{v \times m} s_t + \underbrace{W_0^o}_{v \times 1}) \end{aligned}$$

no transpose

v -class
logistic
regression

Can express as a state machine



- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} x_t + \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} s_{t-1}$$

$$p_t = g(s_t) = f_2(W^o s_t + W_0^o)$$

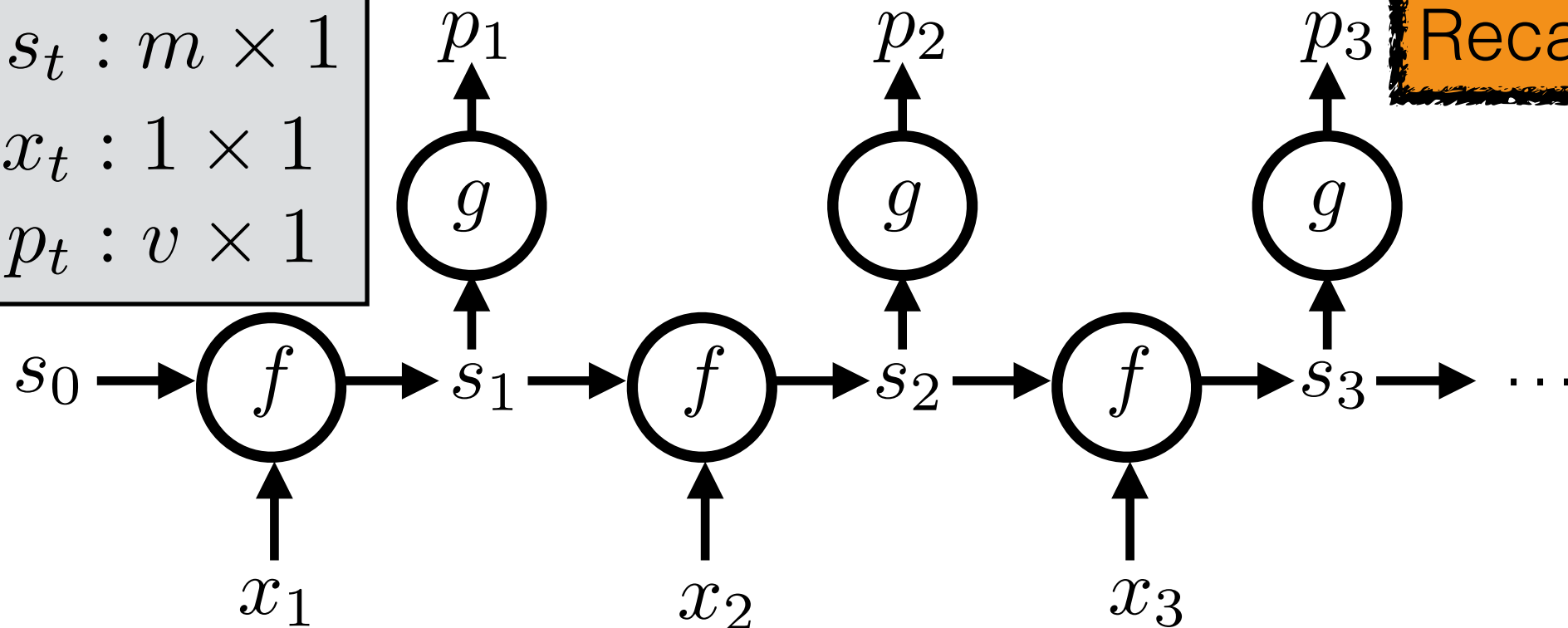
$v \times m$ $v \times 1$

v -class
logistic
regression

Can express as a state machine

$$\begin{array}{l} s_t : m \times 1 \\ x_t : 1 \times 1 \\ p_t : v \times 1 \end{array}$$

Recall: familiar pattern



- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} x_t + \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} s_{t-1}$$

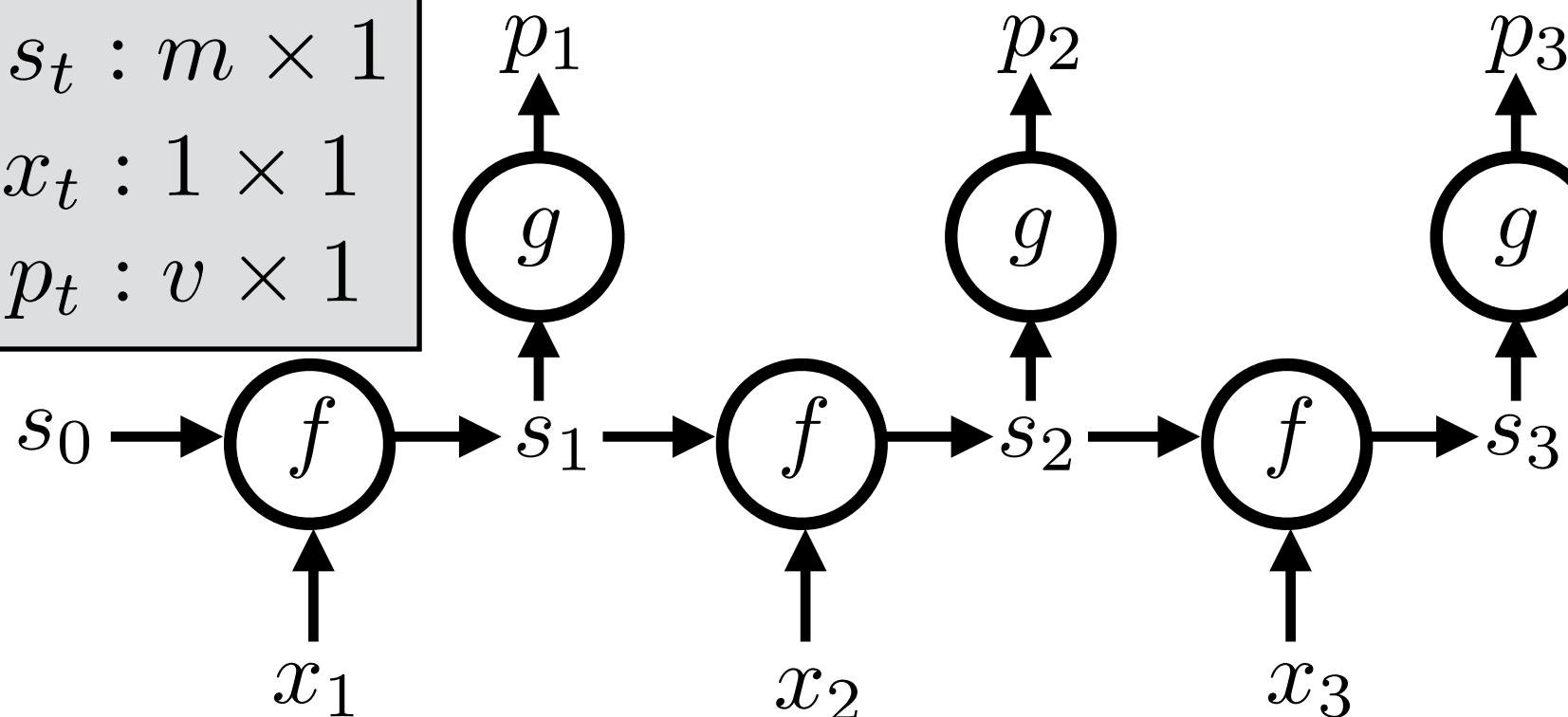
$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(W^o s_t + W_0^o) \end{aligned}$$

$v \times m \quad v \times 1$

v -class
logistic
regression

Can express as a state machine

$$\begin{aligned} s_t &: m \times 1 \\ x_t &: 1 \times 1 \\ p_t &: v \times 1 \end{aligned}$$



Recall: familiar pattern
1. Choose how to predict label (given features & parameters)

- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} x_t + \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} s_{t-1}$$

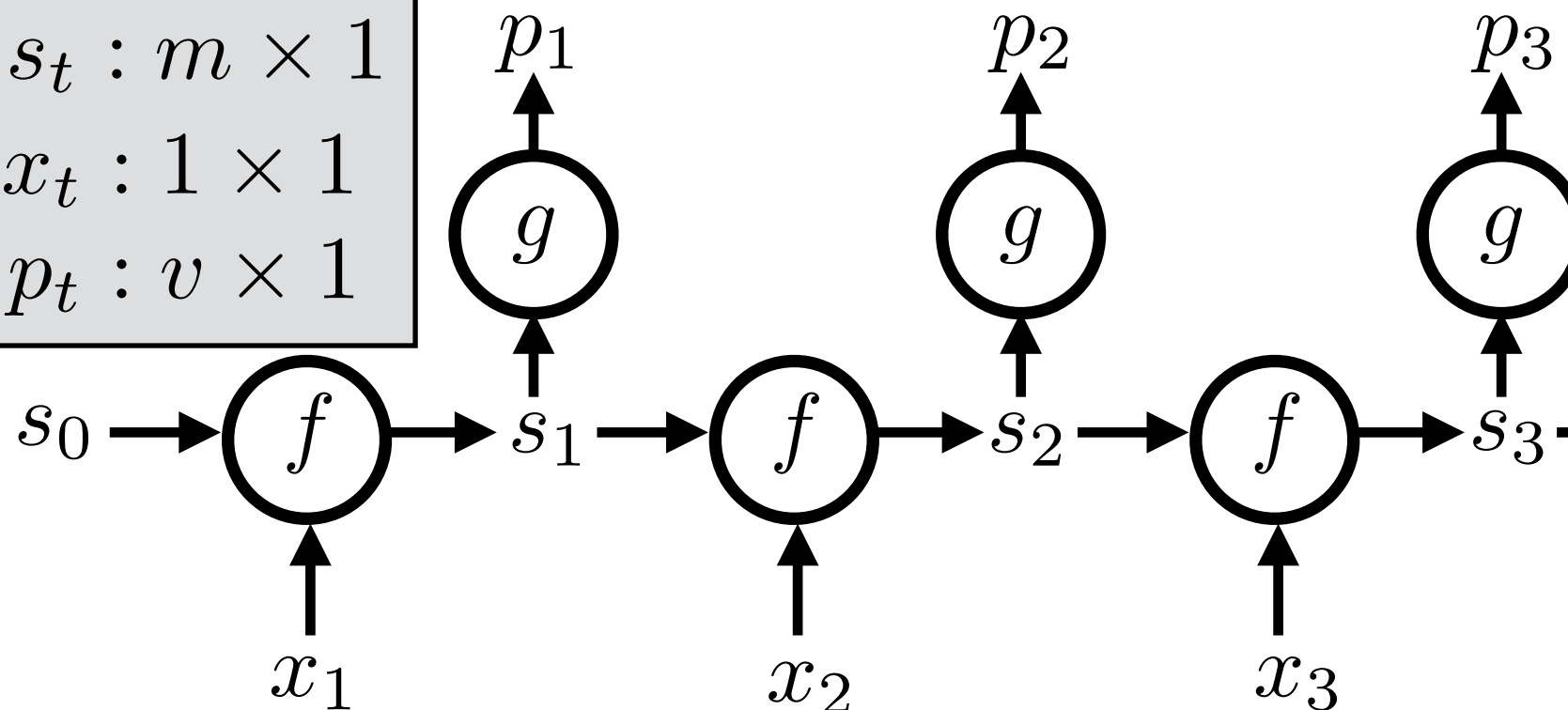
$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(W^o s_t + W_0^o) \end{aligned}$$

$v \times m \quad v \times 1$

v -class
logistic
regression

Can express as a state machine

$$\begin{array}{l} s_t : m \times 1 \\ x_t : 1 \times 1 \\ p_t : v \times 1 \end{array}$$



- Recall: familiar pattern
1. Choose how to predict label (given features & parameters)
 2. Choose a loss (between guess & actual label)

- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} x_t + \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} s_{t-1}$$

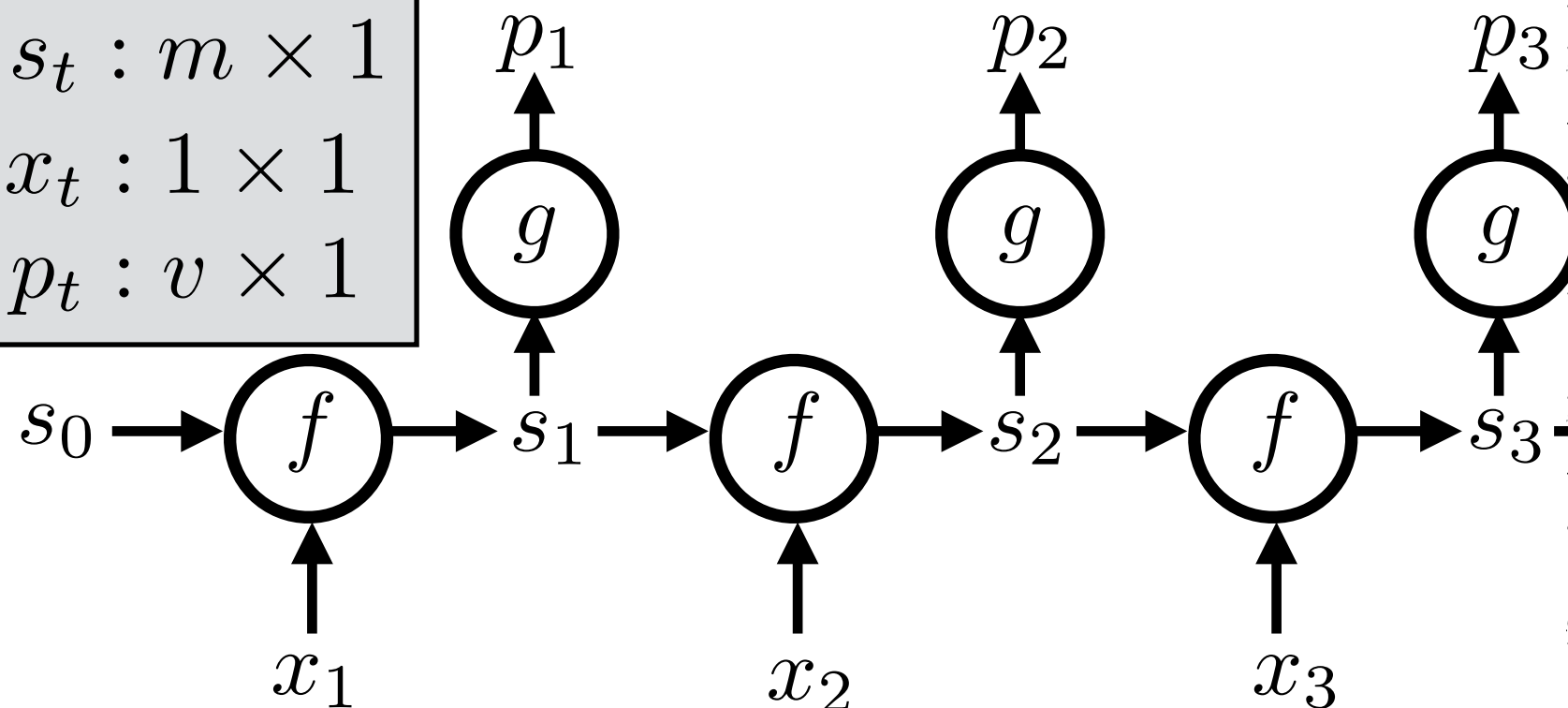
$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(W^o s_t + W_0^o) \end{aligned}$$

$v \times m \quad v \times 1$

v -class
logistic
regression

Can express as a state machine

$$\begin{array}{l} s_t : m \times 1 \\ x_t : 1 \times 1 \\ p_t : v \times 1 \end{array}$$



- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} x_t + \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} s_{t-1}$$

$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(W^o s_t + W_0^o) \end{aligned}$$

$v \times m \quad v \times 1$

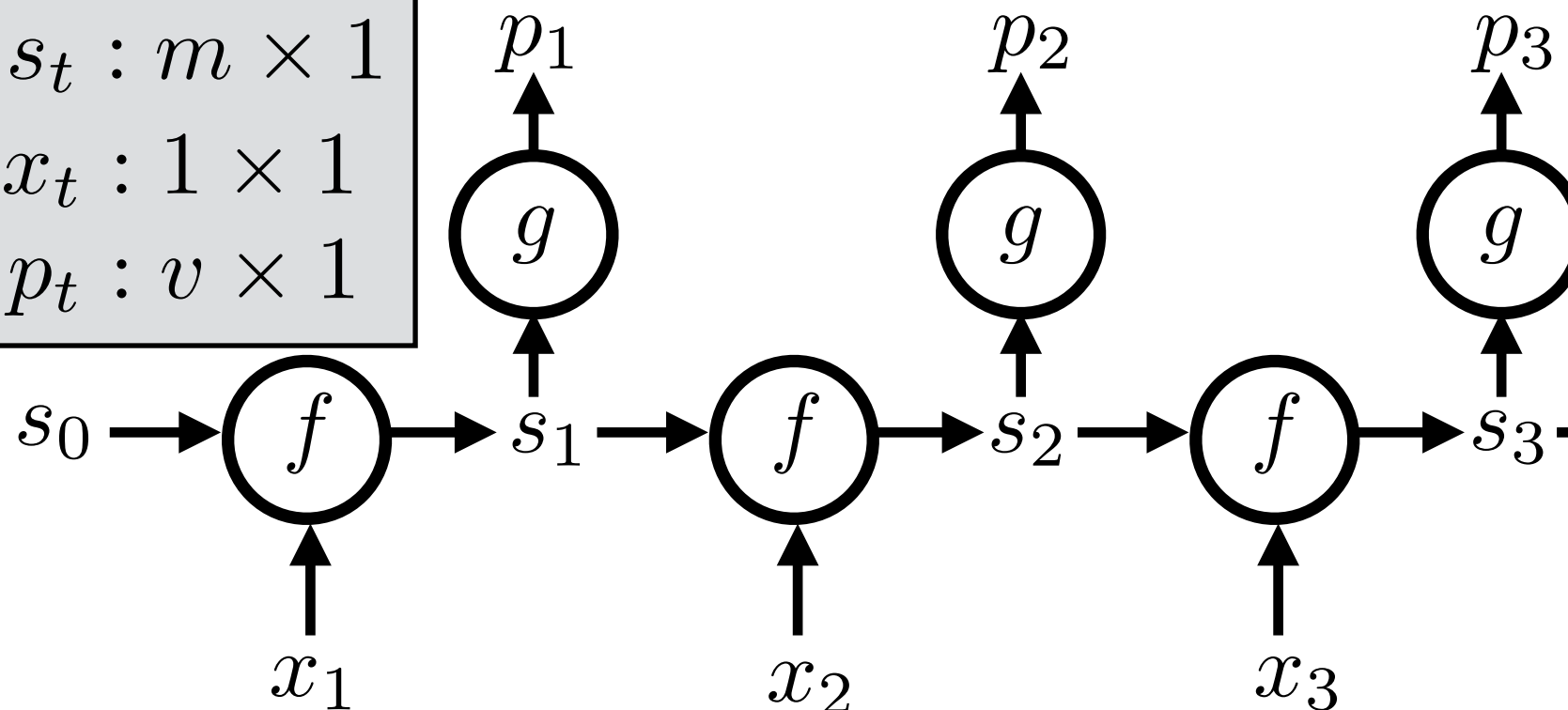
v -class
logistic
regression

Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

Can express as a state machine

$$\begin{array}{l} s_t : m \times 1 \\ x_t : 1 \times 1 \\ p_t : v \times 1 \end{array}$$



- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} x_t + \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} s_{t-1}$$

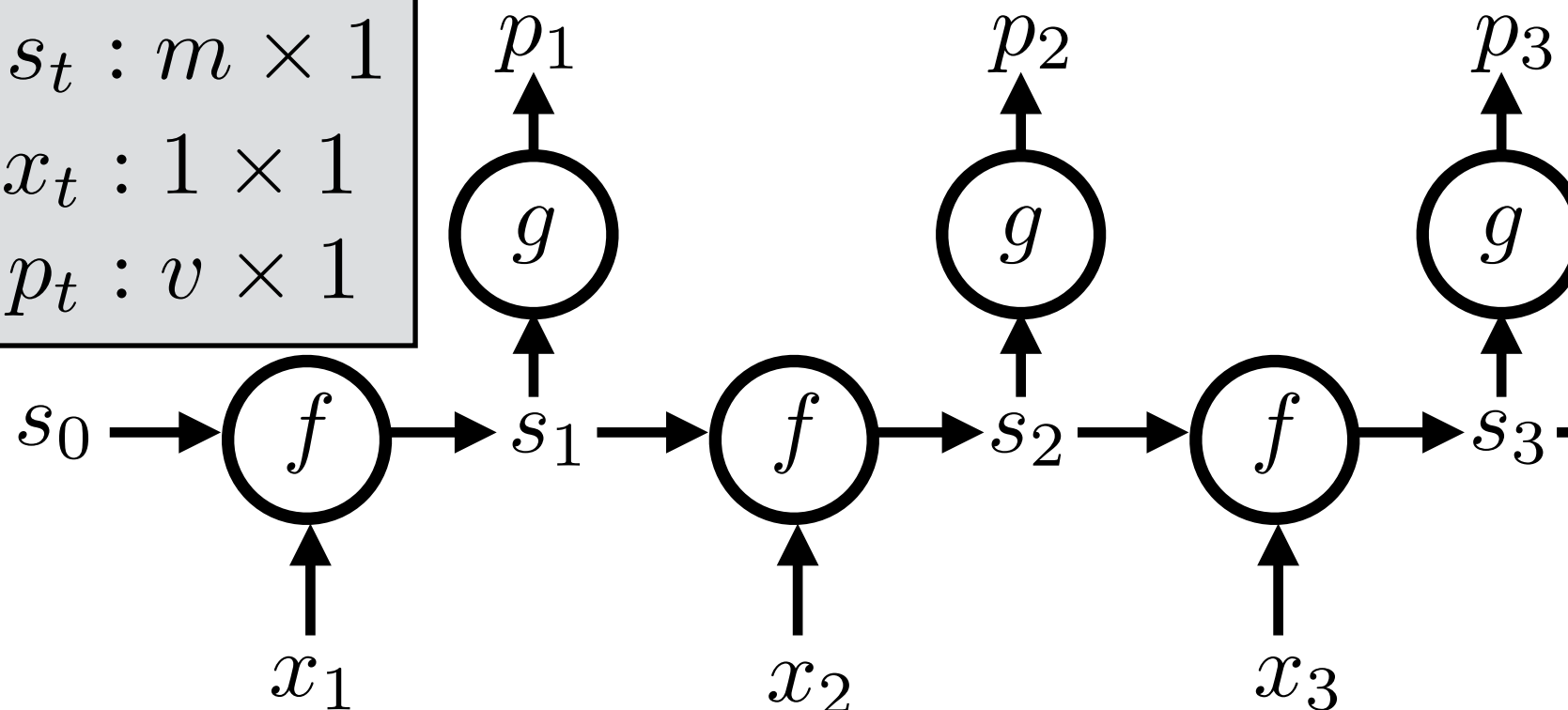
$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(\underbrace{W^o}_{v \times m} s_t + \underbrace{W_0^o}_{v \times 1}) \end{aligned}$$

Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

Can express as a state machine

$$\begin{aligned} s_t &: m \times 1 \\ x_t &: 1 \times 1 \\ p_t &: v \times 1 \end{aligned}$$



- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} x_t + \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} s_{t-1}$$

$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(\underbrace{W^o}_{v \times m} s_t + \underbrace{W_0^o}_{v \times 1}) \end{aligned}$$

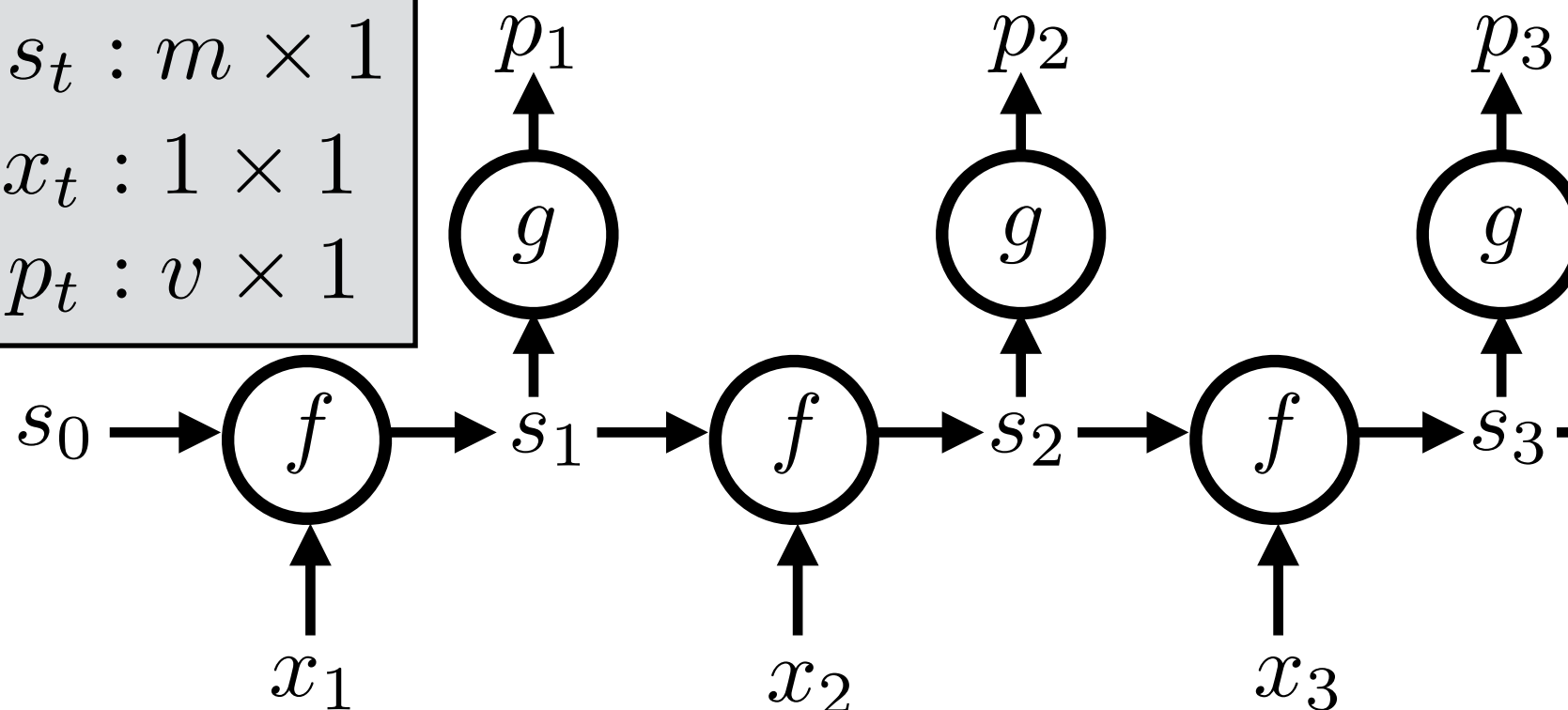
Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

$$L_{\text{elt}}(p_t^{(i)}, y_t^{(i)})$$

Can express as a state machine

$$\begin{aligned} s_t &: m \times 1 \\ x_t &: 1 \times 1 \\ p_t &: v \times 1 \end{aligned}$$



- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} x_t + \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} s_{t-1}$$

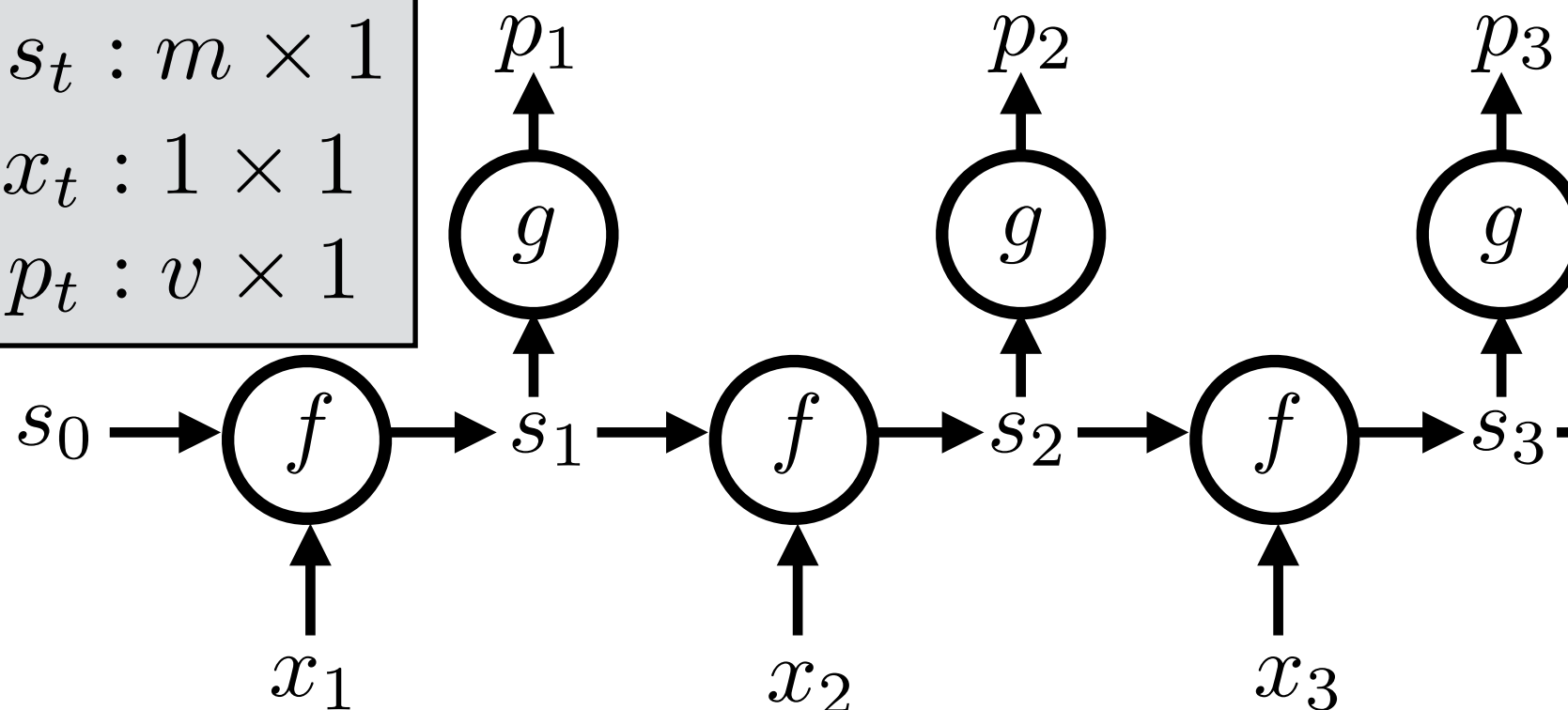
$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(\underbrace{W^o}_{v \times m} s_t + \underbrace{W_0^o}_{v \times 1}) \end{aligned}$$

$$p^{(i)} = R(x^{(i)}; W^o, W_0^o) \quad L_{\text{elt}}(p_t^{(i)}, y_t^{(i)})$$

- Recall: familiar pattern
1. Choose how to predict label (given features & parameters)
 2. Choose a loss (between guess & actual label)
 3. Choose parameters by trying to minimize the training loss

Can express as a state machine

$$\begin{aligned} s_t &: m \times 1 \\ x_t &: 1 \times 1 \\ p_t &: v \times 1 \end{aligned}$$



- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} x_t + \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} s_{t-1}$$

$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(W^o s_t + W_0^o) \end{aligned}$$

$v \times m \quad v \times 1$

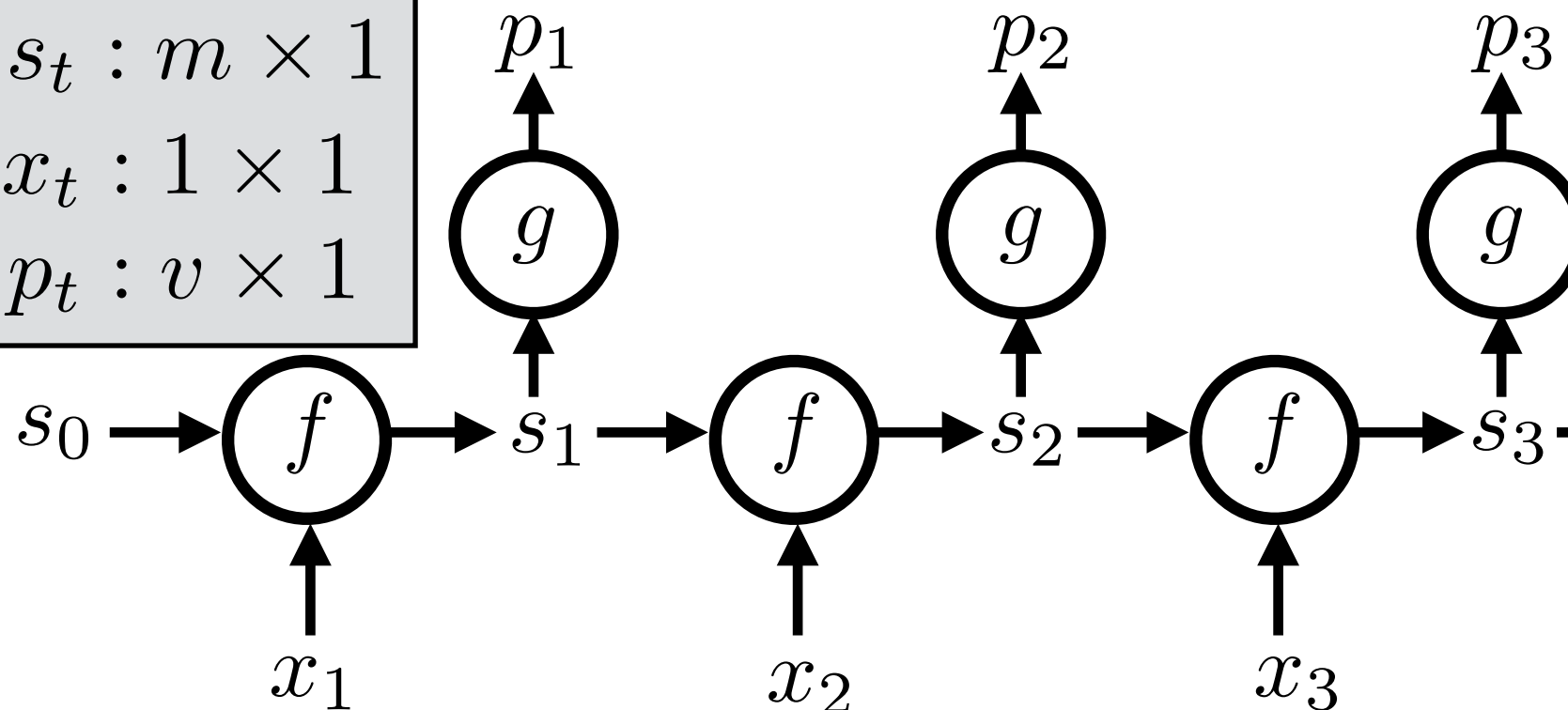
$$\begin{aligned} L_{\text{seq}}(p^{(i)}, y^{(i)}) &= \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)}) \\ p^{(i)} &= R(x^{(i)}; W^o, W_0^o) \end{aligned}$$

Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

Can express as a state machine

$$\begin{aligned} s_t &: m \times 1 \\ x_t &: 1 \times 1 \\ p_t &: v \times 1 \end{aligned}$$



- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} x_t + \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} s_{t-1}$$

$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(W^o s_t + W_0^o) \end{aligned}$$

$v \times m \quad v \times 1$

Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

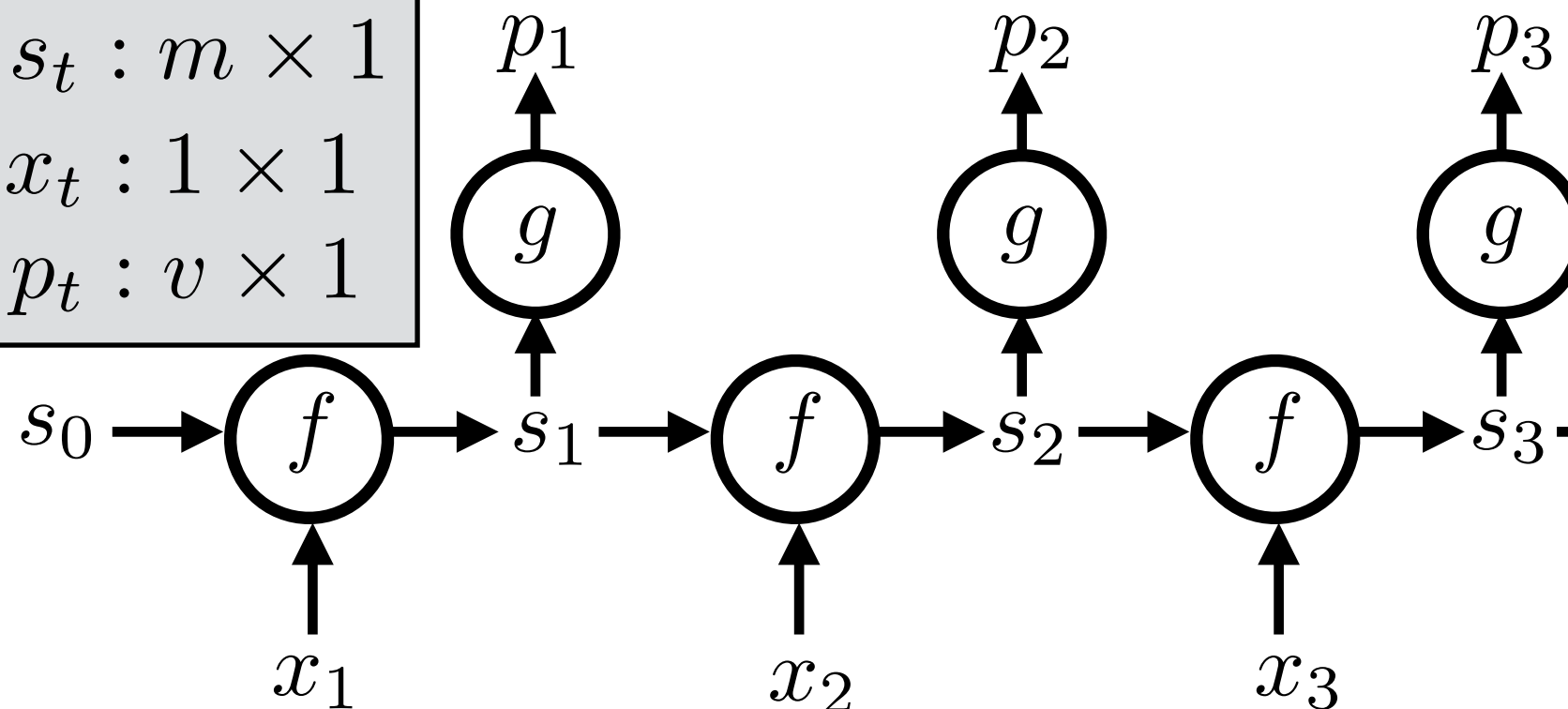
$$L_{\text{seq}}(p^{(i)}, y^{(i)}) = \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)})$$

$$p^{(i)} = R(x^{(i)}; W^o, W_0^o)$$

$$J(W^o, W_0^o) = \sum_{i=1}^q L_{\text{seq}}(p^{(i)}, y^{(i)})$$

Can express as a state machine

$$\begin{aligned} s_t &: m \times 1 \\ x_t &: 1 \times 1 \\ p_t &: v \times 1 \end{aligned}$$



- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} x_t + \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} s_{t-1}$$

$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(W^o s_t + W_0^o) \end{aligned}$$

$v \times m \quad v \times 1$

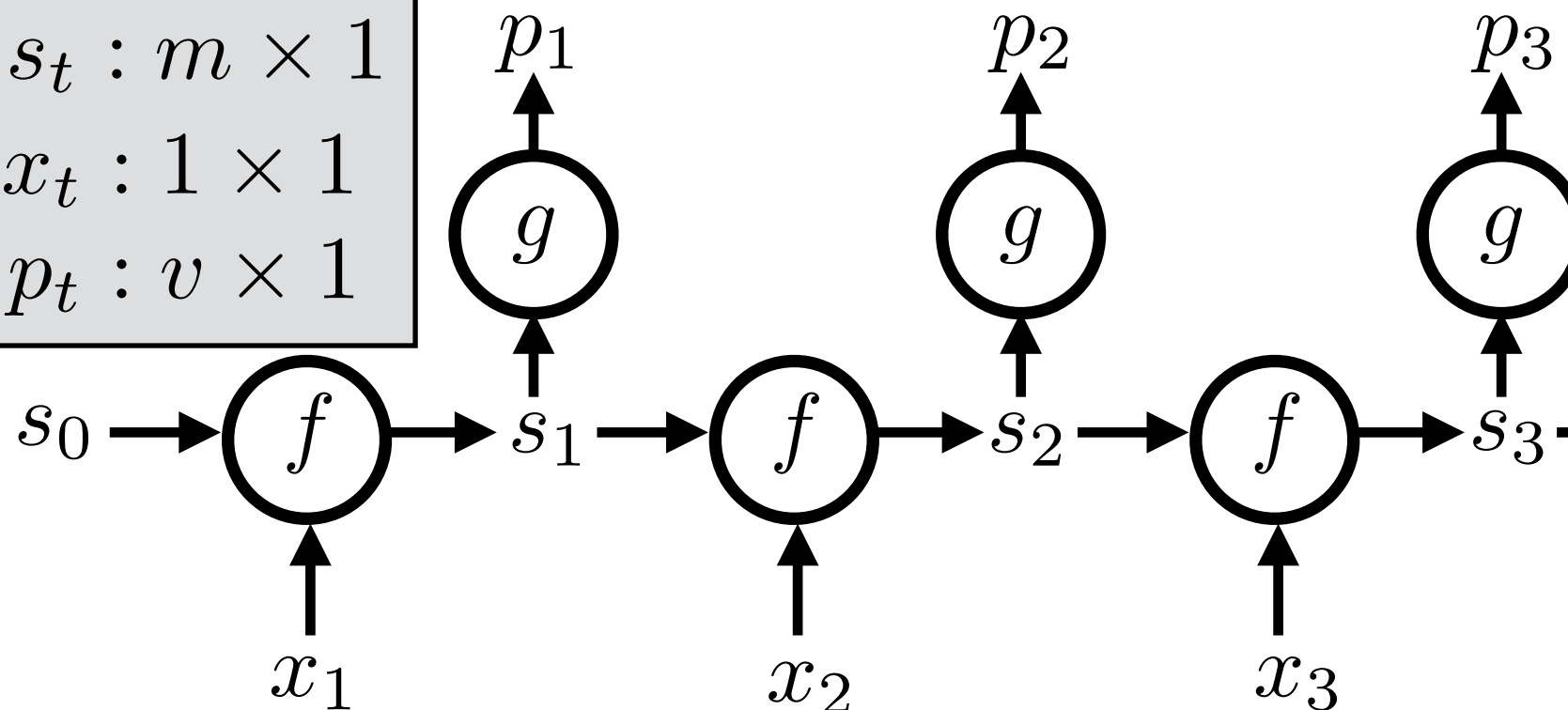
Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

$$\begin{aligned} L_{\text{seq}}(p^{(i)}, y^{(i)}) &= \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)}) \\ p^{(i)} &= R(x^{(i)}; W^o, W_0^o) \\ J(W^o, W_0^o) &= \sum_{i=1}^q L_{\text{seq}}(p^{(i)}, y^{(i)}) \end{aligned}$$

Can express as a state machine

$$\begin{aligned} s_t &: m \times 1 \\ x_t &: 1 \times 1 \\ p_t &: v \times 1 \end{aligned}$$



- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} x_t + \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} s_{t-1}$$

$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(W^o s_t + W_0^o) \end{aligned}$$

$v \times m \quad v \times 1$

Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

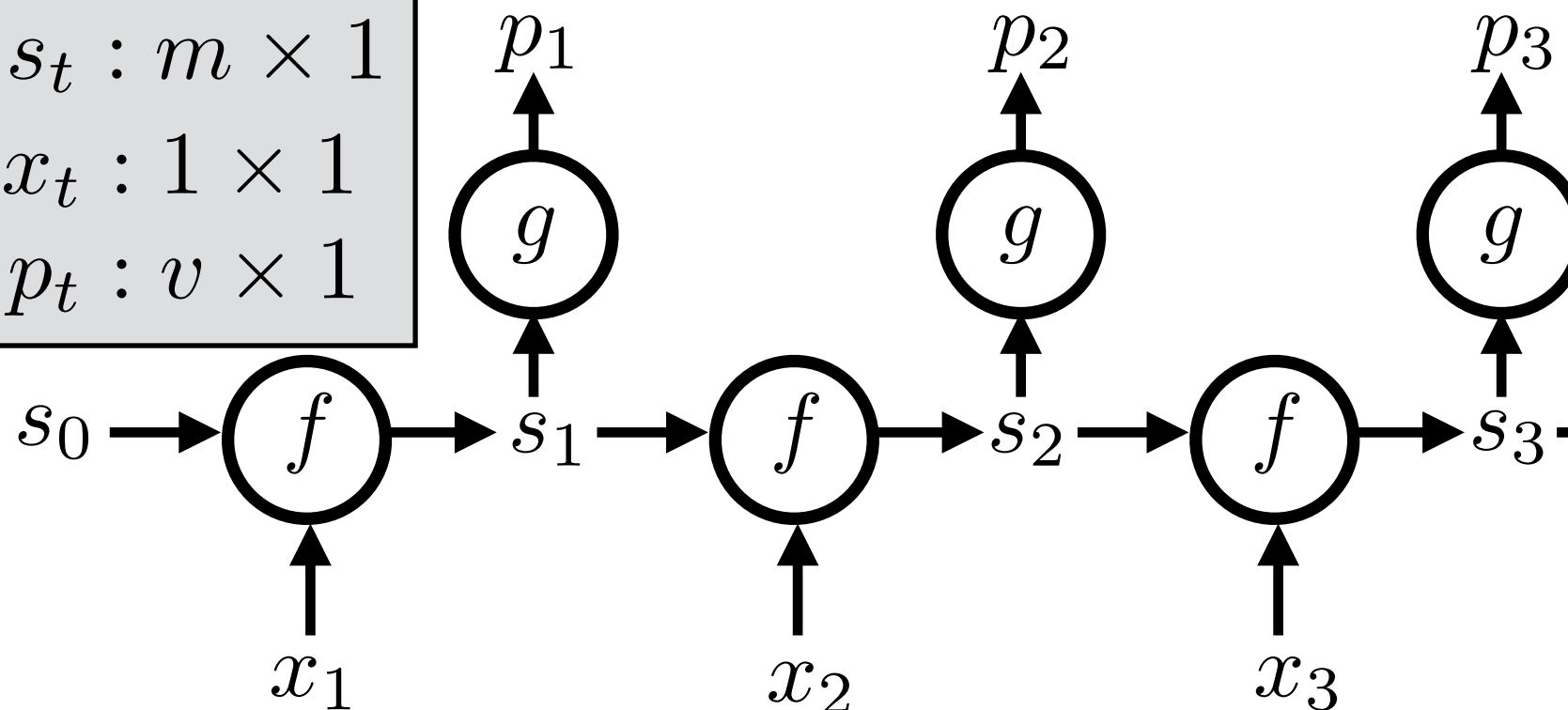
$$L_{\text{seq}}(p^{(i)}, y^{(i)}) = \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)})$$

$$p^{(i)} = R(x^{(i)}; W^o, W_0^o)$$

$$J(W^o, W_0^o) = \sum_{i=1}^q L_{\text{seq}}(p^{(i)}, y^{(i)})$$

Can express as a state machine

$$\begin{aligned} s_t &: m \times 1 \\ x_t &: 1 \times 1 \\ p_t &: v \times 1 \end{aligned}$$



- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} x_t + \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} s_{t-1}$$

$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(W^o s_t + W_0^o) \end{aligned}$$

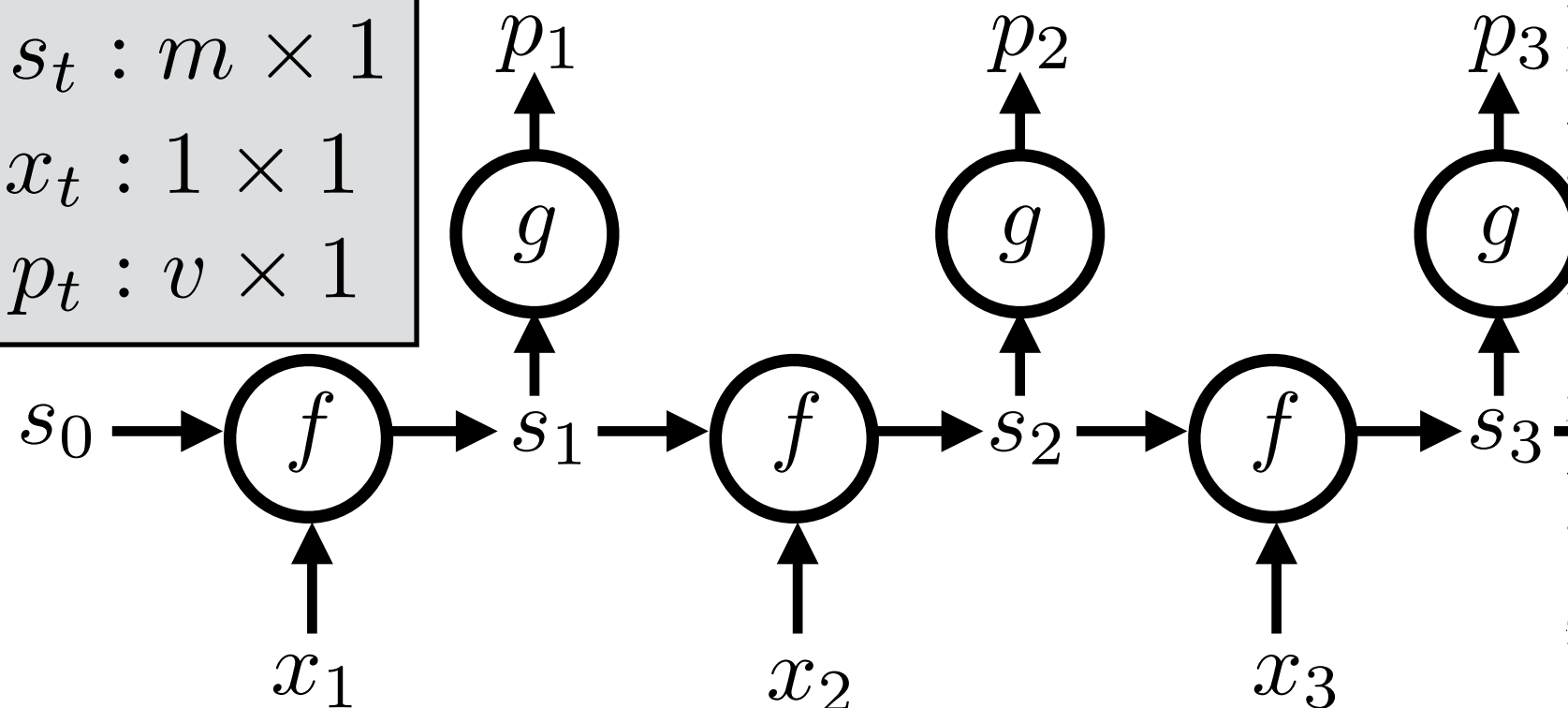
$v \times m \quad v \times 1$

$$\begin{aligned} L_{\text{seq}}(p^{(i)}, y^{(i)}) &= \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)}) \\ p^{(i)} &= R(x^{(i)}; W^o, W_0^o) \\ J(W^o, W_0^o) &= \sum_{i=1}^q L_{\text{seq}}(p^{(i)}, y^{(i)}) \end{aligned}$$

- Recall: familiar pattern
1. Choose how to predict label (given features & parameters)
 2. Choose a loss (between guess & actual label)
 3. Choose parameters by trying to minimize the training loss

Can express as a state machine

$$\begin{aligned} s_t &: m \times 1 \\ x_t &: 1 \times 1 \\ p_t &: v \times 1 \end{aligned}$$



- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} w_1^{sx} \\ w_2^{sx} \\ w_3^{sx} \end{bmatrix} x_t + \begin{bmatrix} w_{11}^{ss} & w_{12}^{ss} & w_{13}^{ss} \\ w_{21}^{ss} & w_{22}^{ss} & w_{23}^{ss} \\ w_{31}^{ss} & w_{32}^{ss} & w_{33}^{ss} \end{bmatrix} s_{t-1}$$

$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(W^o s_t + W_0^o) \end{aligned}$$

$v \times m \quad v \times 1$

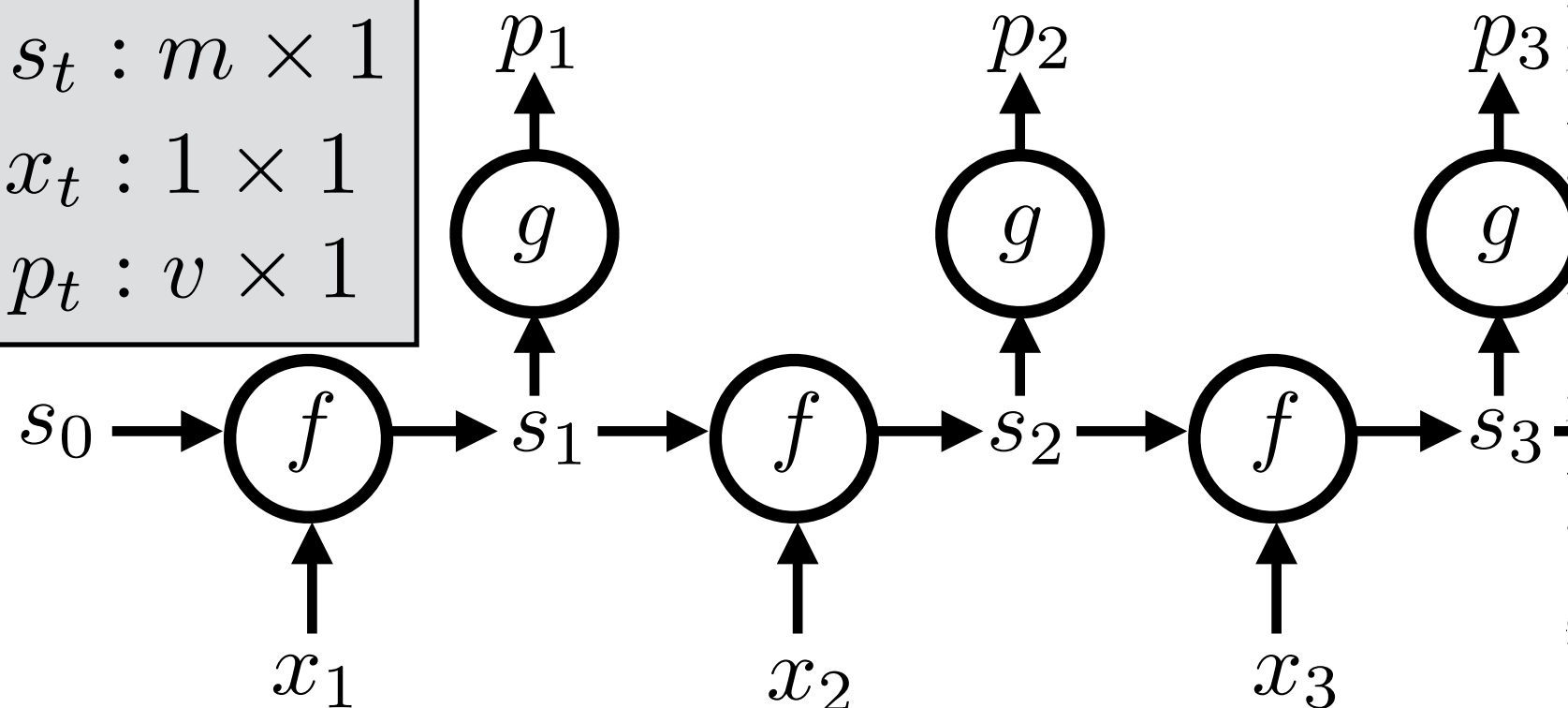
$$L_{\text{seq}}(p^{(i)}, y^{(i)}) = \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)})$$

$$p^{(i)} = R(x^{(i)}; W^o, W_0^o)$$

$$J(W^o, W_0^o) = \sum_{i=1}^q L_{\text{seq}}(p^{(i)}, y^{(i)})$$

Can express as a state machine

$$\begin{aligned} s_t &: m \times 1 \\ x_t &: 1 \times 1 \\ p_t &: v \times 1 \end{aligned}$$



- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} w_1^{sx} \\ w_2^{sx} \\ w_3^{sx} \end{bmatrix} x_t + \begin{bmatrix} w_{11}^{ss} & w_{12}^{ss} & w_{13}^{ss} \\ w_{21}^{ss} & w_{22}^{ss} & w_{23}^{ss} \\ w_{31}^{ss} & w_{32}^{ss} & w_{33}^{ss} \end{bmatrix} s_{t-1}$$

$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(W^o s_t + W_0^o) \end{aligned}$$

$v \times m \quad v \times 1$

Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

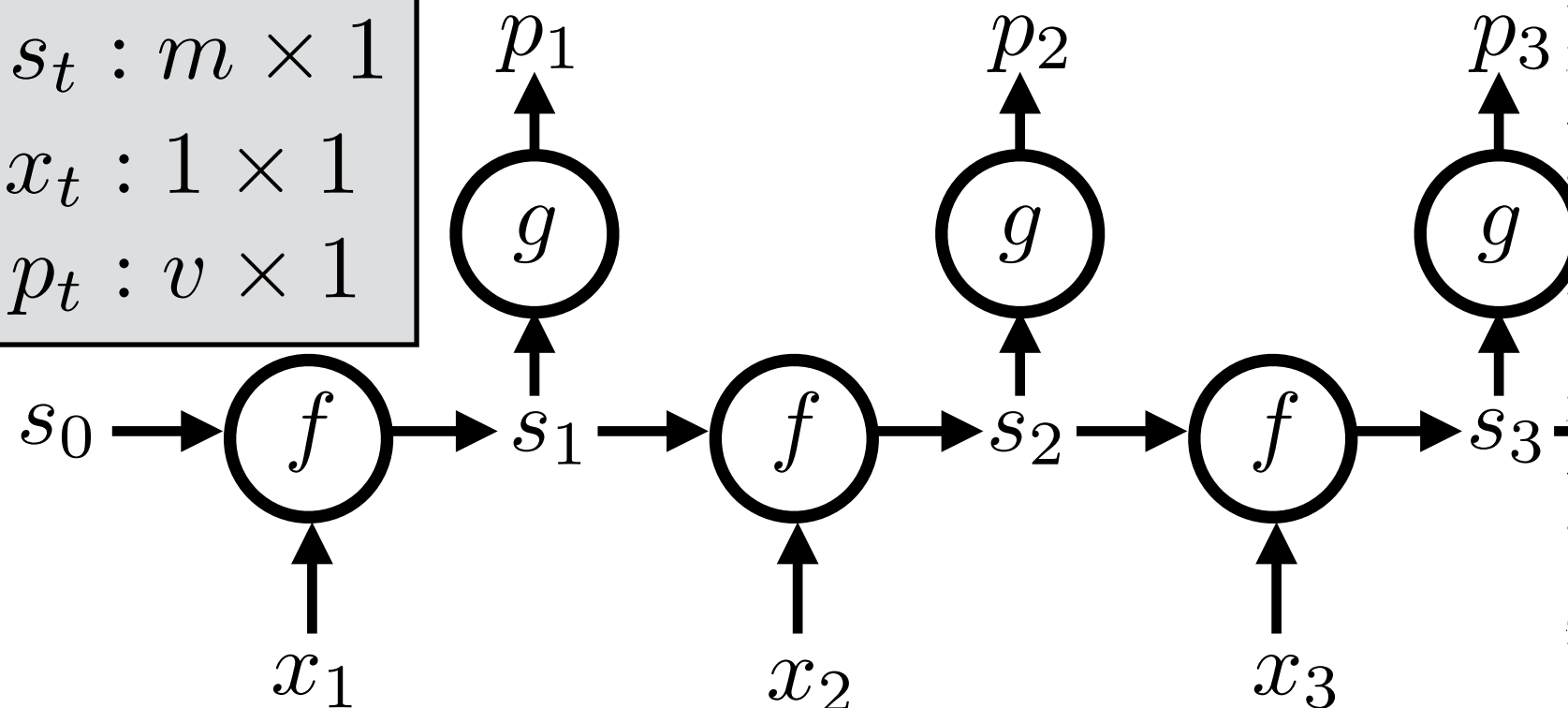
$$L_{\text{seq}}(p^{(i)}, y^{(i)}) = \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)})$$

$$p^{(i)} = R(x^{(i)}; W^o, W_0^o)$$

$$J(W^o, W_0^o) = \sum_{i=1}^q L_{\text{seq}}(p^{(i)}, y^{(i)})$$

Can express as a state machine

$$\begin{aligned} s_t &: m \times 1 \\ x_t &: 1 \times 1 \\ p_t &: v \times 1 \end{aligned}$$



- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} w_1^{sx} \\ w_2^{sx} \\ w_3^{sx} \end{bmatrix} x_t + \begin{bmatrix} w_{11}^{ss} & w_{12}^{ss} & w_{13}^{ss} \\ w_{21}^{ss} & w_{22}^{ss} & w_{23}^{ss} \\ w_{31}^{ss} & w_{32}^{ss} & w_{33}^{ss} \end{bmatrix} s_{t-1}$$

$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(W^o s_t + W_0^o) \end{aligned}$$

$v \times m \quad v \times 1$

Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

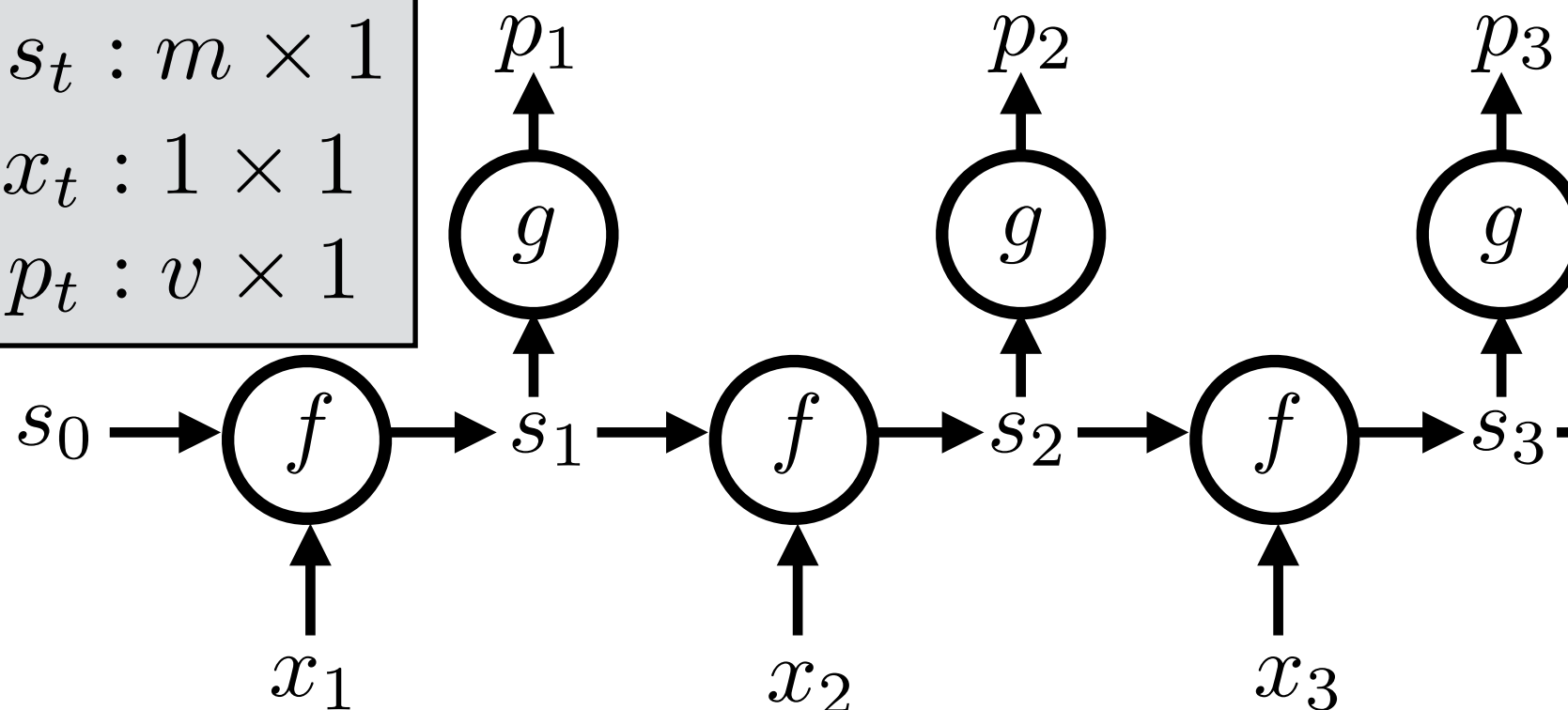
$$L_{\text{seq}}(p^{(i)}, y^{(i)}) = \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)})$$

$$p^{(i)} = R(x^{(i)}; W^o, W_0^o)$$

$$J(W^o, W_0^o) = \sum_{i=1}^q L_{\text{seq}}(p^{(i)}, y^{(i)})$$

Can express as a state machine

$$\begin{aligned} s_t &: m \times 1 \\ x_t &: 1 \times 1 \\ p_t &: v \times 1 \end{aligned}$$



- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} w_1^{sx} \\ w_2^{sx} \\ w_3^{sx} \end{bmatrix} x_t + \begin{bmatrix} w_{11}^{ss} & w_{12}^{ss} & w_{13}^{ss} \\ w_{21}^{ss} & w_{22}^{ss} & w_{23}^{ss} \\ w_{31}^{ss} & w_{32}^{ss} & w_{33}^{ss} \end{bmatrix} s_{t-1}$$

$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(\underbrace{W^o}_{v \times m} s_t + \underbrace{W_0^o}_{v \times 1}) \end{aligned}$$

Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

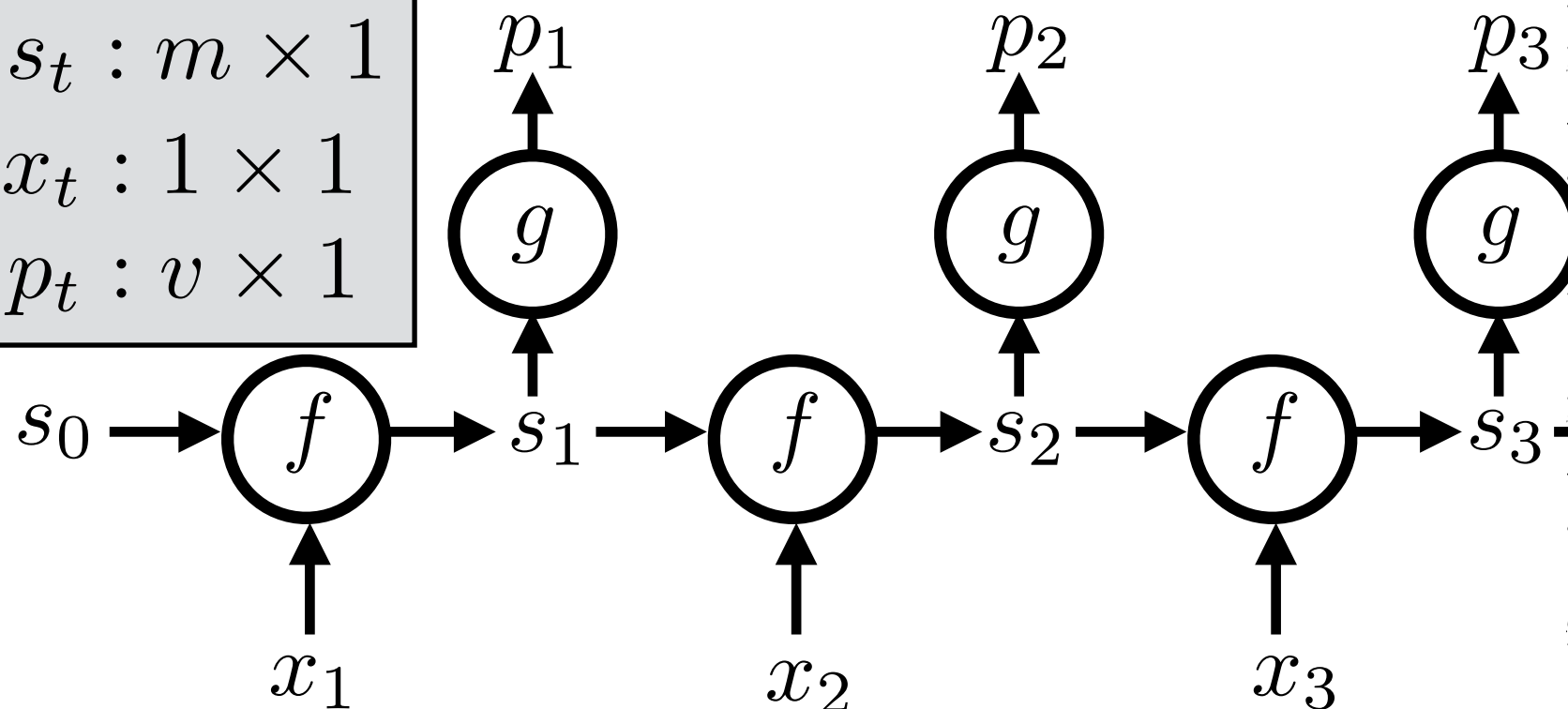
$$L_{\text{seq}}(p^{(i)}, y^{(i)}) = \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)})$$

$$p^{(i)} = R(x^{(i)}; W^o, W_0^o)$$

$$J(W^o, W_0^o) = \sum_{i=1}^q L_{\text{seq}}(p^{(i)}, y^{(i)})$$

Can express as a state machine

$$\begin{aligned} s_t &: m \times 1 \\ x_t &: 1 \times 1 \\ p_t &: v \times 1 \end{aligned}$$



- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = \begin{bmatrix} w_1^{sx} \\ w_2^{sx} \\ w_3^{sx} \end{bmatrix} x_t + \begin{bmatrix} w_{11}^{ss} & w_{12}^{ss} & w_{13}^{ss} \\ w_{21}^{ss} & w_{22}^{ss} & w_{23}^{ss} \\ w_{31}^{ss} & w_{32}^{ss} & w_{33}^{ss} \end{bmatrix} s_{t-1}$$

$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(\underbrace{W^o}_{v \times m} s_t + \underbrace{W_0^o}_{v \times 1}) \end{aligned}$$

Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

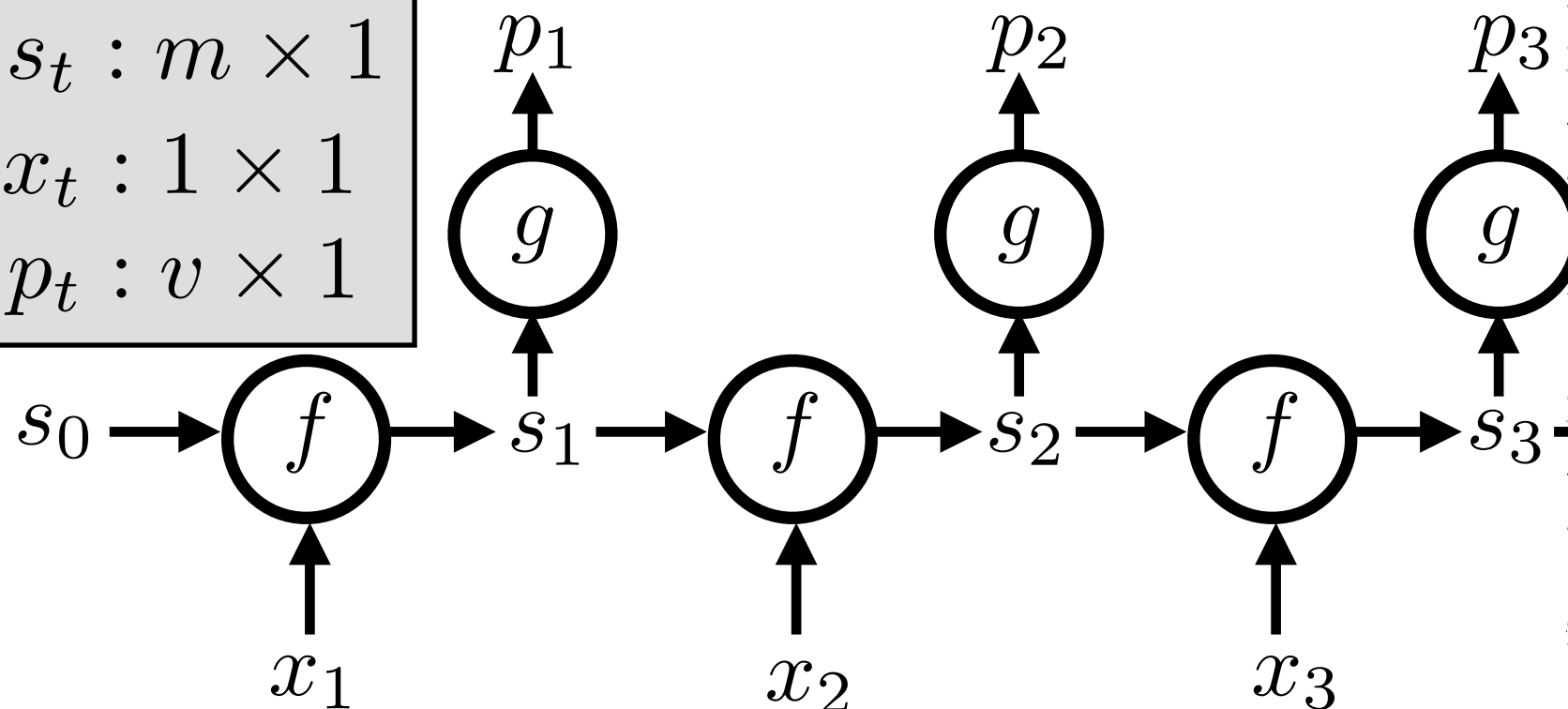
$$L_{\text{seq}}(p^{(i)}, y^{(i)}) = \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)})$$

$$p^{(i)} = R(x^{(i)}; W^o, W_0^o)$$

$$J(W^o, W_0^o) = \sum_{i=1}^q L_{\text{seq}}(p^{(i)}, y^{(i)})$$

Can express as a state machine

$$\begin{aligned} s_t &: m \times 1 \\ x_t &: 1 \times 1 \\ p_t &: v \times 1 \end{aligned}$$



- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = \begin{bmatrix} w_1^{sx} \\ w_2^{sx} \\ w_3^{sx} \end{bmatrix} x_t + \begin{bmatrix} w_{11}^{ss} & w_{12}^{ss} & w_{13}^{ss} \\ w_{21}^{ss} & w_{22}^{ss} & w_{23}^{ss} \\ w_{31}^{ss} & w_{32}^{ss} & w_{33}^{ss} \end{bmatrix} s_{t-1} + \begin{bmatrix} w_{0,1}^{ss} \\ w_{0,2}^{ss} \\ w_{0,3}^{ss} \end{bmatrix}$$

$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(\underbrace{W^o}_{v \times m} s_t + \underbrace{W_0^o}_{v \times 1}) \end{aligned}$$

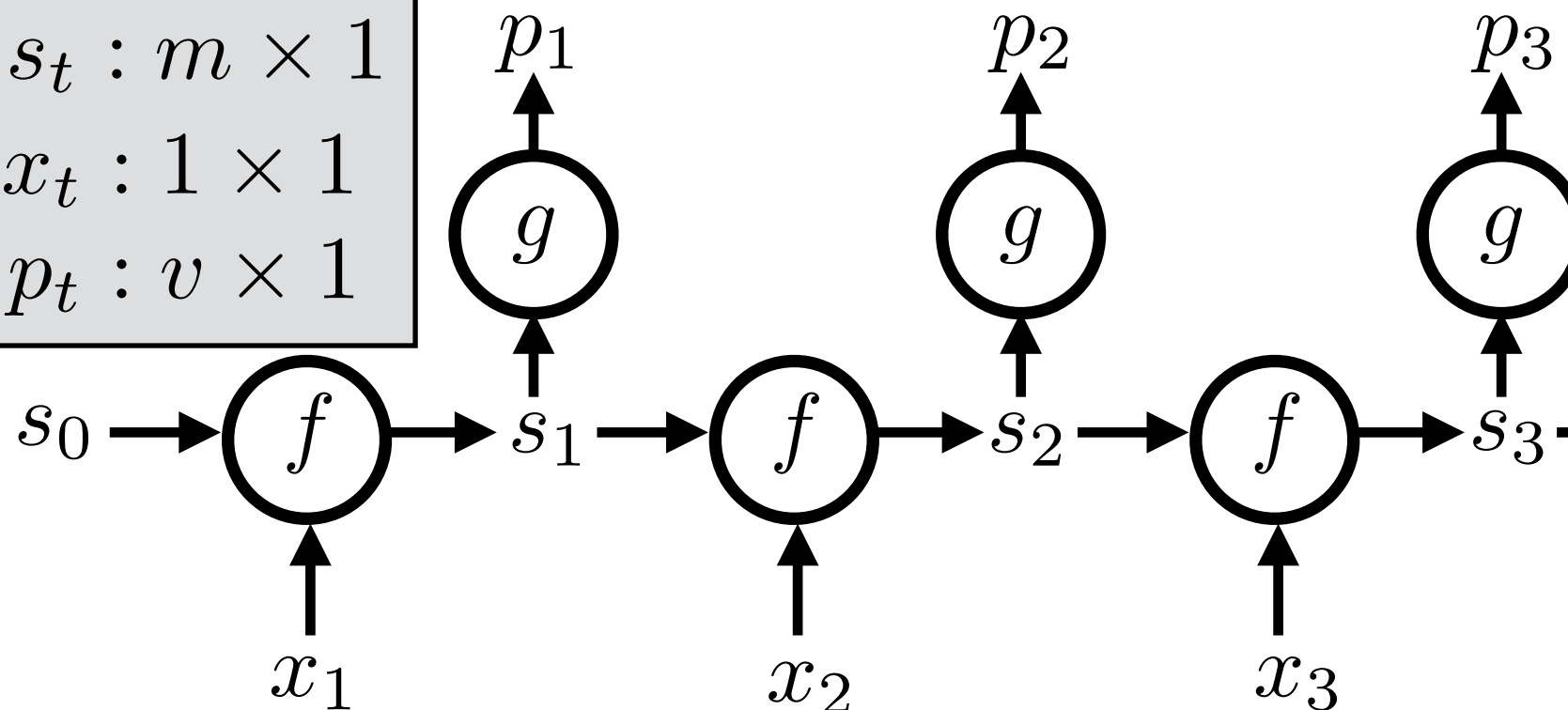
$$\begin{aligned} L_{\text{seq}}(p^{(i)}, y^{(i)}) &= \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)}) \\ p^{(i)} &= R(x^{(i)}; W^o, W_0^o) \\ J(W^o, W_0^o) &= \sum_{i=1}^q L_{\text{seq}}(p^{(i)}, y^{(i)}) \end{aligned}$$

Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

Can express as a state machine

$$\begin{aligned} s_t &: m \times 1 \\ x_t &: 1 \times 1 \\ p_t &: v \times 1 \end{aligned}$$



- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = \begin{bmatrix} w_1^{sx} \\ w_2^{sx} \\ w_3^{sx} \end{bmatrix} x_t + \begin{bmatrix} w_{11}^{ss} & w_{12}^{ss} & w_{13}^{ss} \\ w_{21}^{ss} & w_{22}^{ss} & w_{23}^{ss} \\ w_{31}^{ss} & w_{32}^{ss} & w_{33}^{ss} \end{bmatrix} s_{t-1} + \begin{bmatrix} w_{0,1}^{ss} \\ w_{0,2}^{ss} \\ w_{0,3}^{ss} \end{bmatrix}$$

$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(\underbrace{W^o}_{v \times m} s_t + \underbrace{W_0^o}_{v \times 1}) \end{aligned}$$

Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

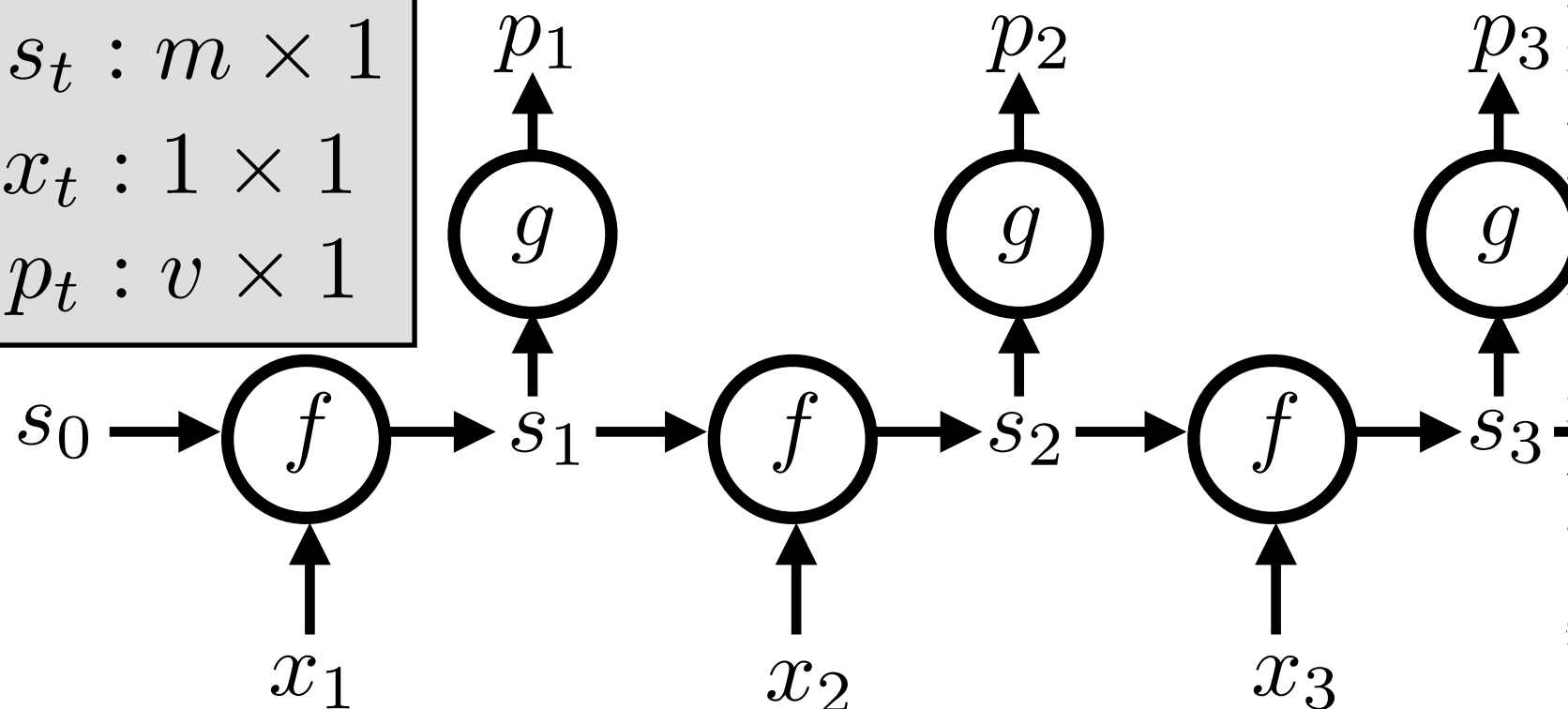
$$L_{\text{seq}}(p^{(i)}, y^{(i)}) = \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)})$$

$$p^{(i)} = R(x^{(i)}; W^o, W_0^o)$$

$$J(W^o, W_0^o) = \sum_{i=1}^q L_{\text{seq}}(p^{(i)}, y^{(i)})$$

Can express as a state machine

$$\begin{aligned} s_t &: m \times 1 \\ x_t &: 1 \times 1 \\ p_t &: v \times 1 \end{aligned}$$



- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = \begin{bmatrix} w_1^{sx} \\ w_2^{sx} \\ w_3^{sx} \end{bmatrix} x_t + \begin{bmatrix} w_{11}^{ss} & w_{12}^{ss} & w_{13}^{ss} \\ w_{21}^{ss} & w_{22}^{ss} & w_{23}^{ss} \\ w_{31}^{ss} & w_{32}^{ss} & w_{33}^{ss} \end{bmatrix} s_{t-1} + \begin{bmatrix} w_{0,1}^{ss} \\ w_{0,2}^{ss} \\ w_{0,3}^{ss} \end{bmatrix}$$

$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(W^o s_t + W_0^o) \end{aligned}$$

$v \times m \quad v \times 1$

Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

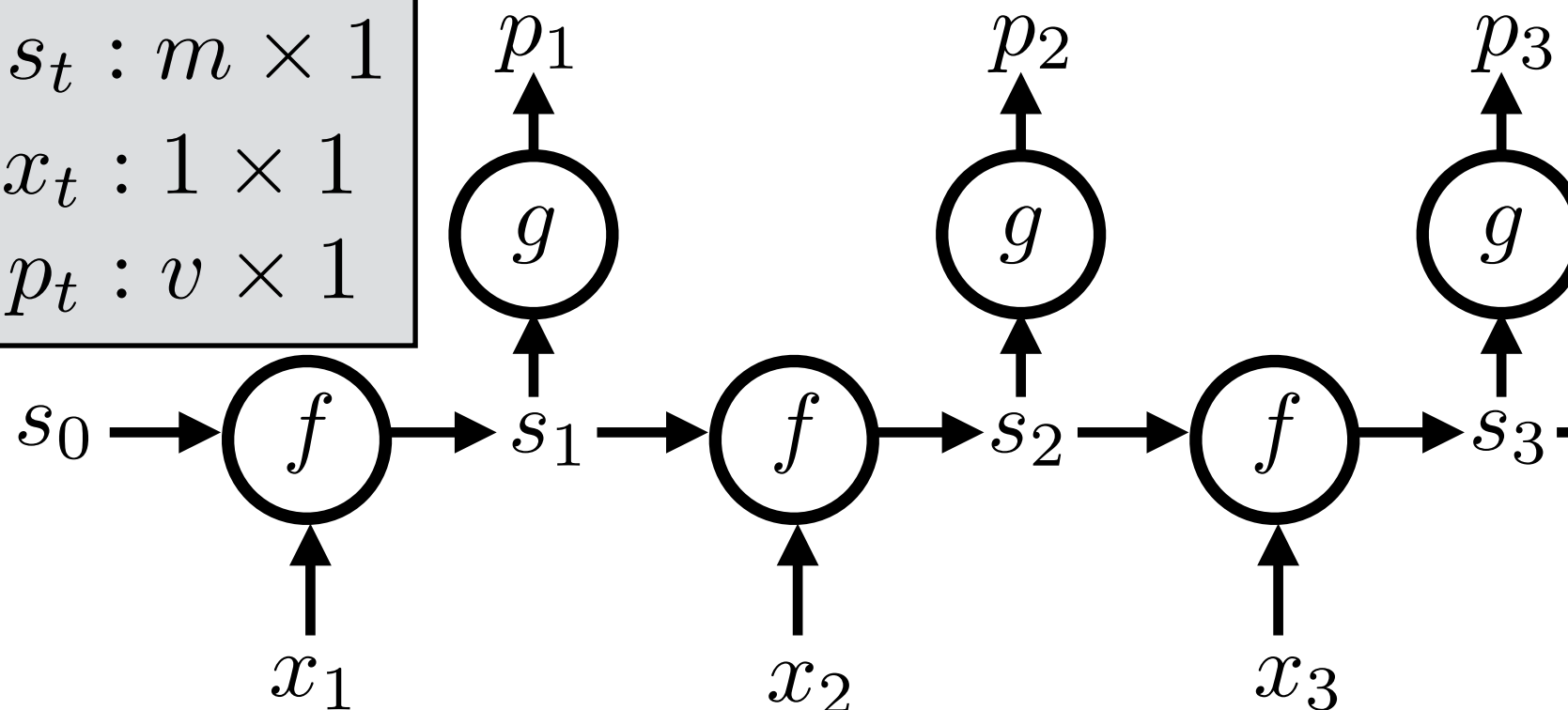
$$L_{\text{seq}}(p^{(i)}, y^{(i)}) = \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)})$$

$$p^{(i)} = R(x^{(i)}; W^o, W_0^o)$$

$$J(W^o, W_0^o) = \sum_{i=1}^q L_{\text{seq}}(p^{(i)}, y^{(i)})$$

Can express as a state machine

$$\begin{aligned} s_t &: m \times 1 \\ x_t &: 1 \times 1 \\ p_t &: v \times 1 \end{aligned}$$



- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = \begin{bmatrix} w_1^{sx} \\ w_2^{sx} \\ w_3^{sx} \end{bmatrix} x_t + \begin{bmatrix} w_{11}^{ss} & w_{12}^{ss} & w_{13}^{ss} \\ w_{21}^{ss} & w_{22}^{ss} & w_{23}^{ss} \\ w_{31}^{ss} & w_{32}^{ss} & w_{33}^{ss} \end{bmatrix} s_{t-1} + \begin{bmatrix} w_{0,1}^{ss} \\ w_{0,2}^{ss} \\ w_{0,3}^{ss} \end{bmatrix}$$

$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(W^o s_t + W_0^o) \end{aligned}$$

$v \times m \quad v \times 1$

Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

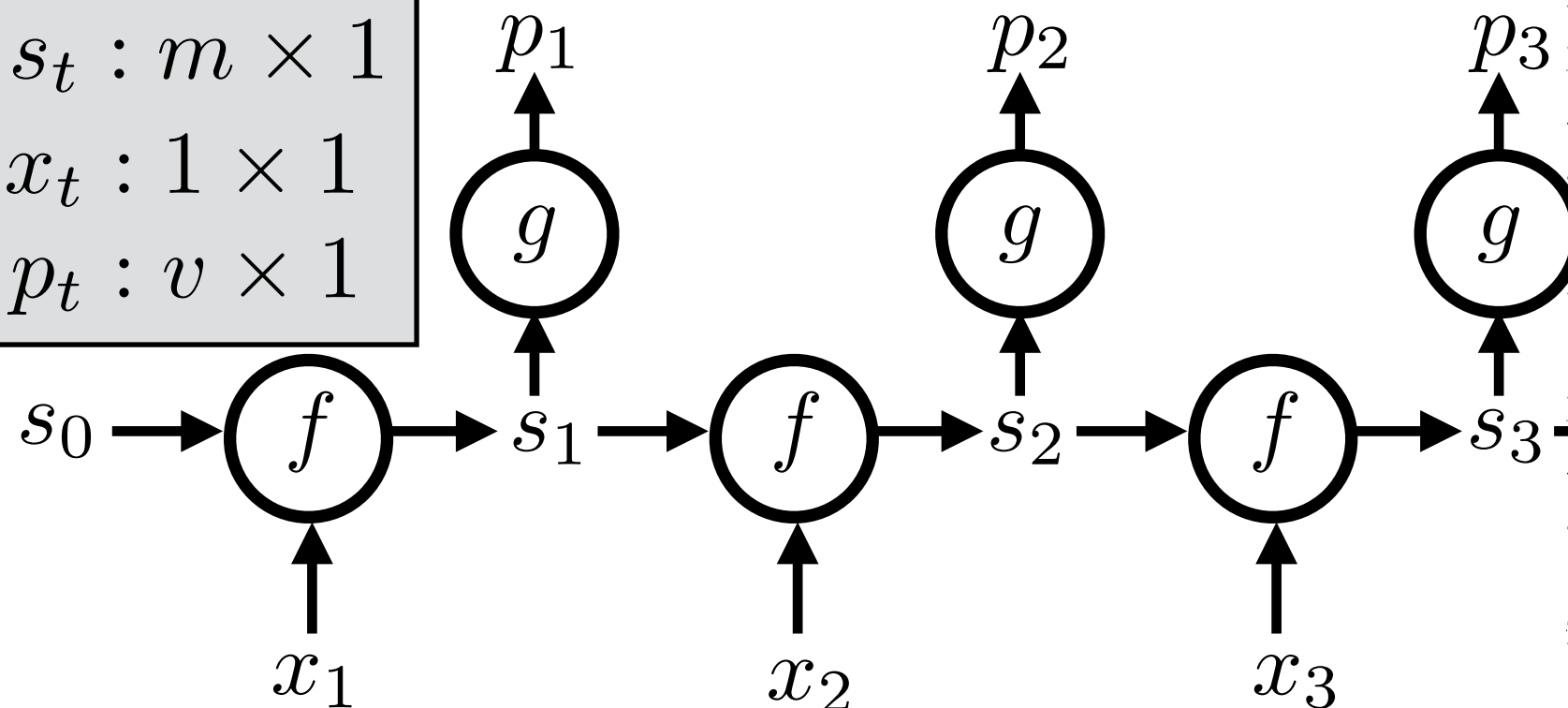
$$L_{\text{seq}}(p^{(i)}, y^{(i)}) = \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)})$$

$$p^{(i)} = R(x^{(i)}; W^o, W_0^o)$$

$$J(W^o, W_0^o) = \sum_{i=1}^q L_{\text{seq}}(p^{(i)}, y^{(i)})$$

Can express as a state machine

$$\begin{aligned} s_t &: m \times 1 \\ x_t &: 1 \times 1 \\ p_t &: v \times 1 \end{aligned}$$



- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = f_1 \left(\begin{bmatrix} w_1^{sx} \\ w_2^{sx} \\ w_3^{sx} \end{bmatrix} x_t + \begin{bmatrix} w_{11}^{ss} & w_{12}^{ss} & w_{13}^{ss} \\ w_{21}^{ss} & w_{22}^{ss} & w_{23}^{ss} \\ w_{31}^{ss} & w_{32}^{ss} & w_{33}^{ss} \end{bmatrix} s_{t-1} + \begin{bmatrix} w_{0,1}^{ss} \\ w_{0,2}^{ss} \\ w_{0,3}^{ss} \end{bmatrix} \right)$$

$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(W^o s_t + W_0^o) \\ &\quad \quad \quad v \times m \quad \quad v \times 1 \end{aligned}$$

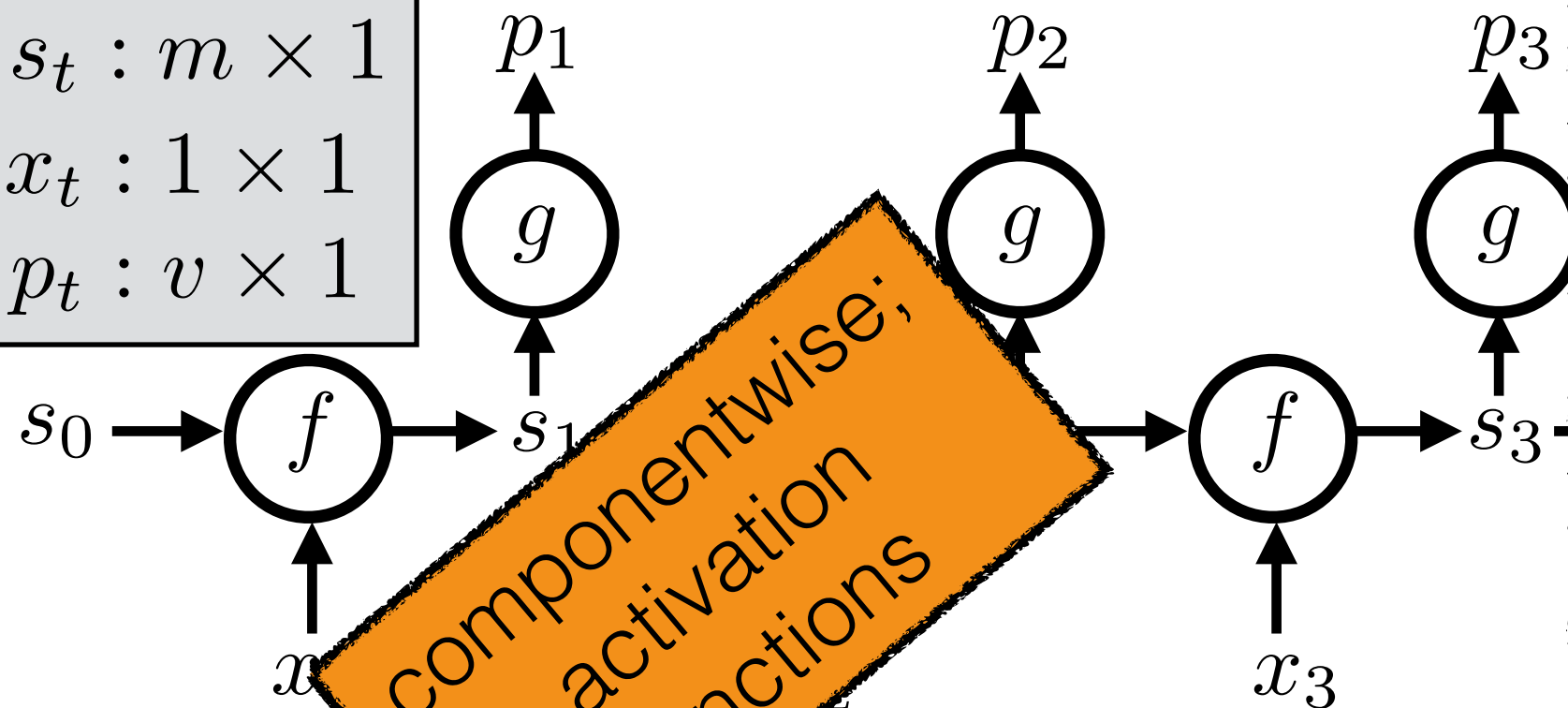
Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

$$\begin{aligned} L_{\text{seq}}(p^{(i)}, y^{(i)}) &= \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)}) \\ p^{(i)} &= R(x^{(i)}; W^o, W_0^o) \\ J(W^o, W_0^o) &= \sum_{i=1}^q L_{\text{seq}}(p^{(i)}, y^{(i)}) \end{aligned}$$

Can express as a state machine

$$\begin{aligned} s_t &: m \times 1 \\ x_t &: 1 \times 1 \\ p_t &: v \times 1 \end{aligned}$$



- Example: state is last set $\{0, 1\}$; $v = 3$ characters

$$s_t = f_1 \left(\begin{bmatrix} w_1^{sx} \\ w_2^{sx} \\ w_3^{sx} \end{bmatrix} x_t + \begin{bmatrix} w_{11}^{ss} & w_{12}^{ss} & w_{13}^{ss} \\ w_{21}^{ss} & w_{22}^{ss} & w_{23}^{ss} \\ w_{31}^{ss} & w_{32}^{ss} & w_{33}^{ss} \end{bmatrix} s_{t-1} + \begin{bmatrix} w_{0,1}^{ss} \\ w_{0,2}^{ss} \\ w_{0,3}^{ss} \end{bmatrix} \right)$$

$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(W^o s_t + W_0^o) \end{aligned}$$

$v \times m \quad v \times 1$

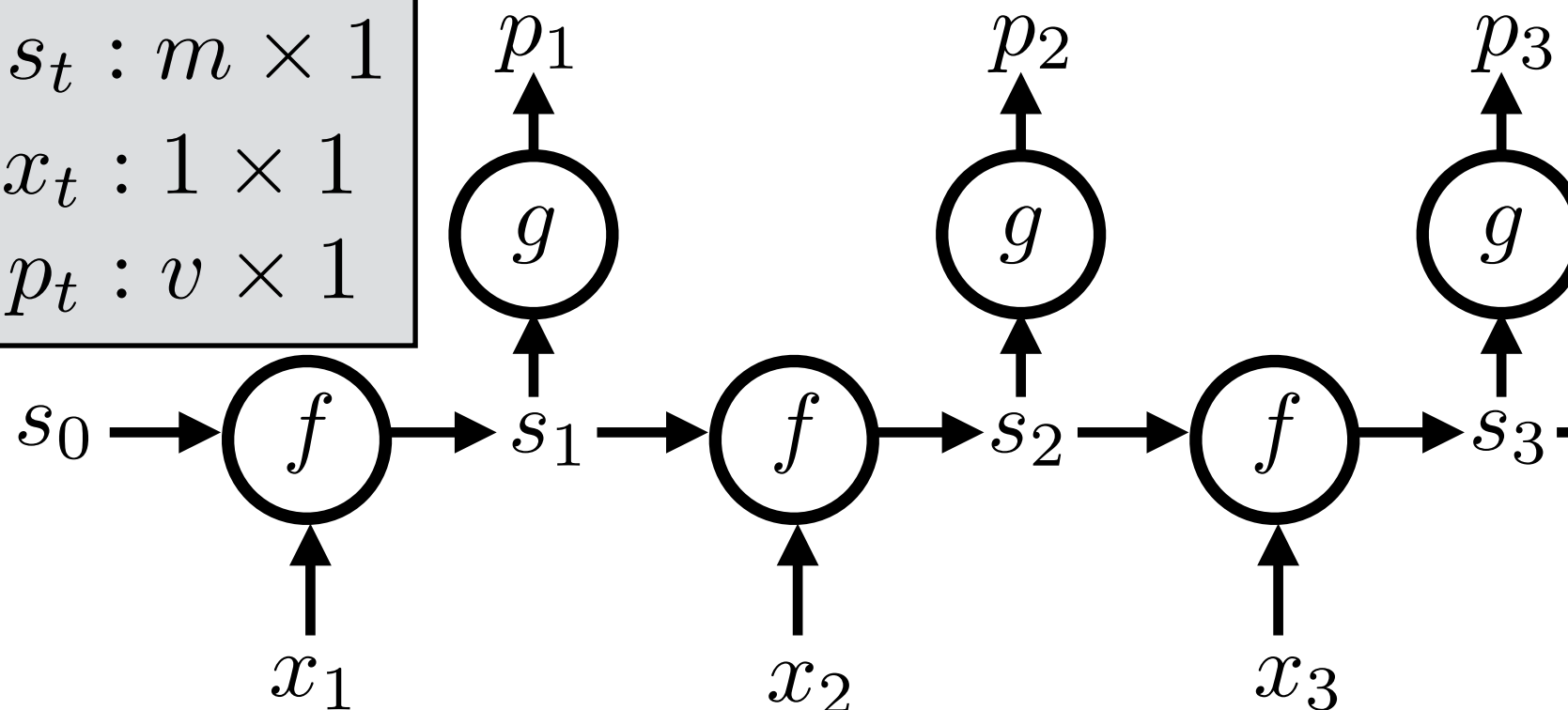
Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

$$\begin{aligned} L_{\text{seq}}(p^{(i)}, y^{(i)}) &= \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)}) \\ p^{(i)} &= R(x^{(i)}; W^o, W_0^o) \\ J(W^o, W_0^o) &= \sum_{i=1}^q L_{\text{seq}}(p^{(i)}, y^{(i)}) \end{aligned}$$

Can express as a state machine

$$\begin{aligned} s_t &: m \times 1 \\ x_t &: 1 \times 1 \\ p_t &: v \times 1 \end{aligned}$$



- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = f_1 \left(\begin{bmatrix} w_1^{sx} \\ w_2^{sx} \\ w_3^{sx} \end{bmatrix} x_t + \begin{bmatrix} w_{11}^{ss} & w_{12}^{ss} & w_{13}^{ss} \\ w_{21}^{ss} & w_{22}^{ss} & w_{23}^{ss} \\ w_{31}^{ss} & w_{32}^{ss} & w_{33}^{ss} \end{bmatrix} s_{t-1} + \begin{bmatrix} w_{0,1}^{ss} \\ w_{0,2}^{ss} \\ w_{0,3}^{ss} \end{bmatrix} \right)$$

$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(W^o s_t + W_0^o) \\ &\quad \quad \quad v \times m \quad \quad v \times 1 \end{aligned}$$

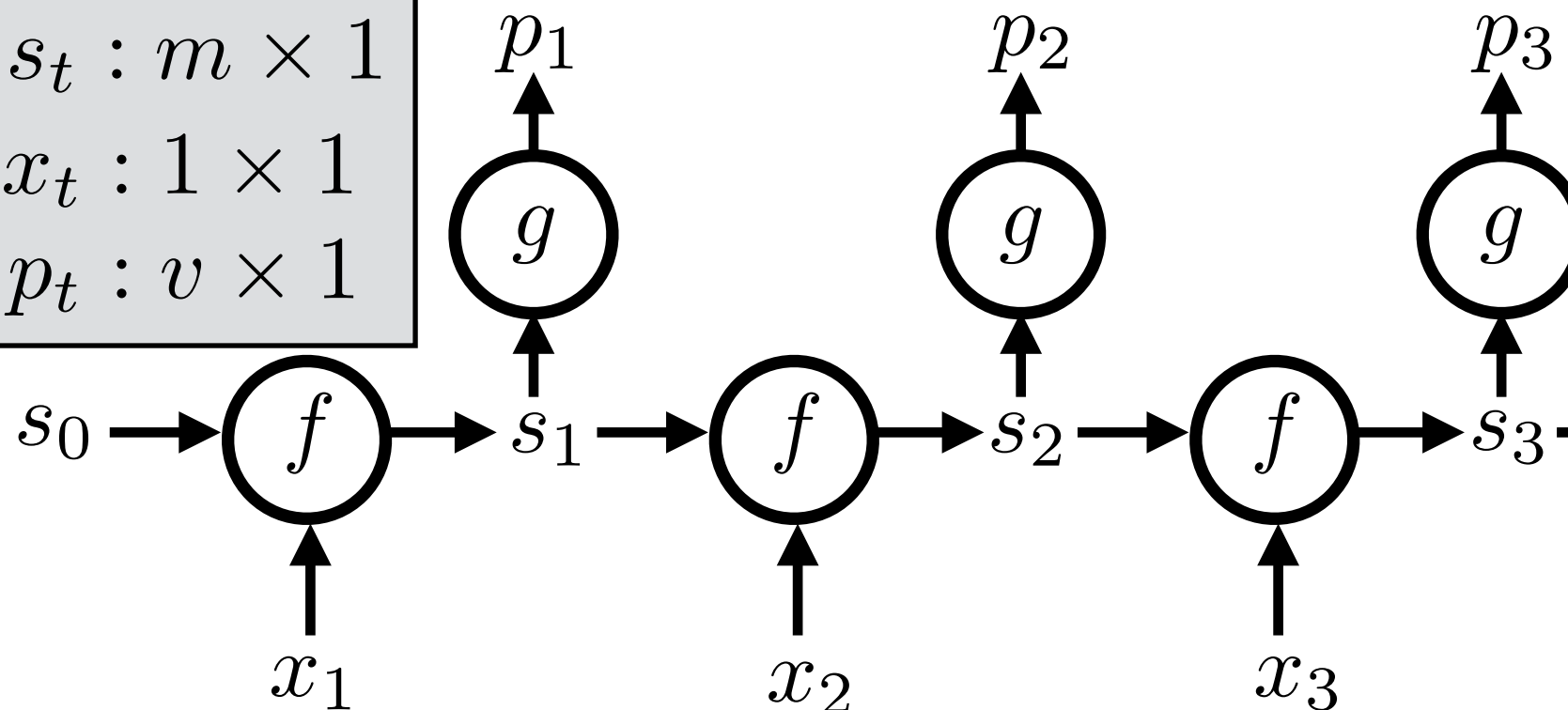
Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

$$\begin{aligned} L_{\text{seq}}(p^{(i)}, y^{(i)}) &= \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)}) \\ p^{(i)} &= R(x^{(i)}; W^o, W_0^o) \\ J(W^o, W_0^o) &= \sum_{i=1}^q L_{\text{seq}}(p^{(i)}, y^{(i)}) \end{aligned}$$

Can express as a state machine

$$\begin{aligned} s_t &: m \times 1 \\ x_t &: 1 \times 1 \\ p_t &: v \times 1 \end{aligned}$$



- Example: Alphabet $\{0,1\}$;
state is last $m = 3$ characters

$$s_t = f_1 \left(\begin{bmatrix} w_1^{sx} \\ w_2^{sx} \\ w_3^{sx} \end{bmatrix} x_t + \begin{bmatrix} w_{11}^{ss} & w_{12}^{ss} & w_{13}^{ss} \\ w_{21}^{ss} & w_{22}^{ss} & w_{23}^{ss} \\ w_{31}^{ss} & w_{32}^{ss} & w_{33}^{ss} \end{bmatrix} s_{t-1} + \begin{bmatrix} w_{0,1}^{ss} \\ w_{0,2}^{ss} \\ w_{0,3}^{ss} \end{bmatrix} \right)$$

$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(\underbrace{W^o}_{v \times m} s_t + \underbrace{W_0^o}_{v \times 1}) \end{aligned}$$

Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

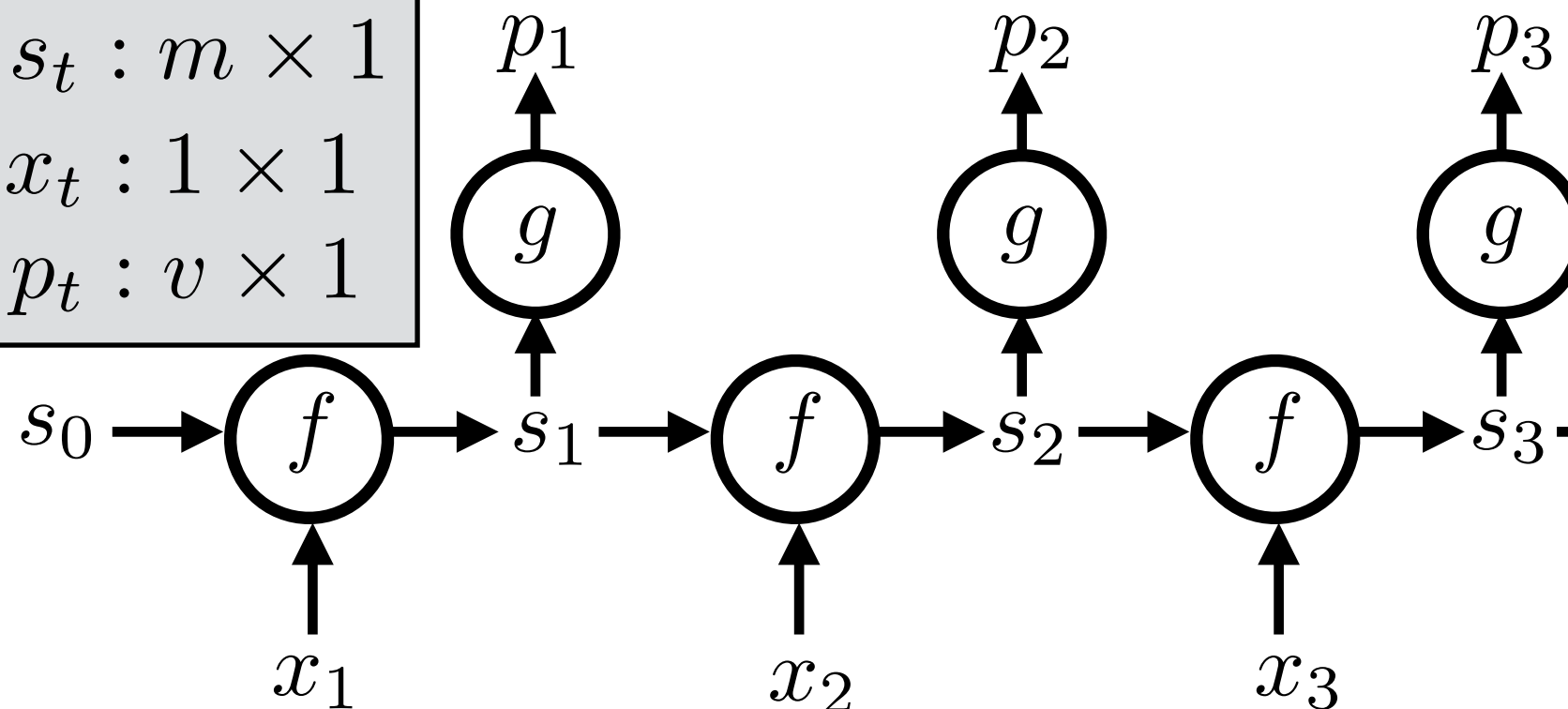
$$L_{\text{seq}}(p^{(i)}, y^{(i)}) = \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)})$$

$$p^{(i)} = R(x^{(i)}; W^o, W_0^o)$$

$$J(W^o, W_0^o) = \sum_{i=1}^q L_{\text{seq}}(p^{(i)}, y^{(i)})$$

Can express as a state machine

$$\begin{aligned} s_t &: m \times 1 \\ x_t &: 1 \times 1 \\ p_t &: v \times 1 \end{aligned}$$



- Example: Alphabet $\{0, 1\}$;
state is last $m = 3$ characters

- Recall: familiar pattern
1. Choose how to predict label (given features & parameters)
 2. Choose a loss (between guess & actual label)
 3. Choose parameters by trying to minimize the training loss

$$s_t = f_1 (W^{sx} x_t + W^{ss} s_{t-1} + W_0^{ss})$$

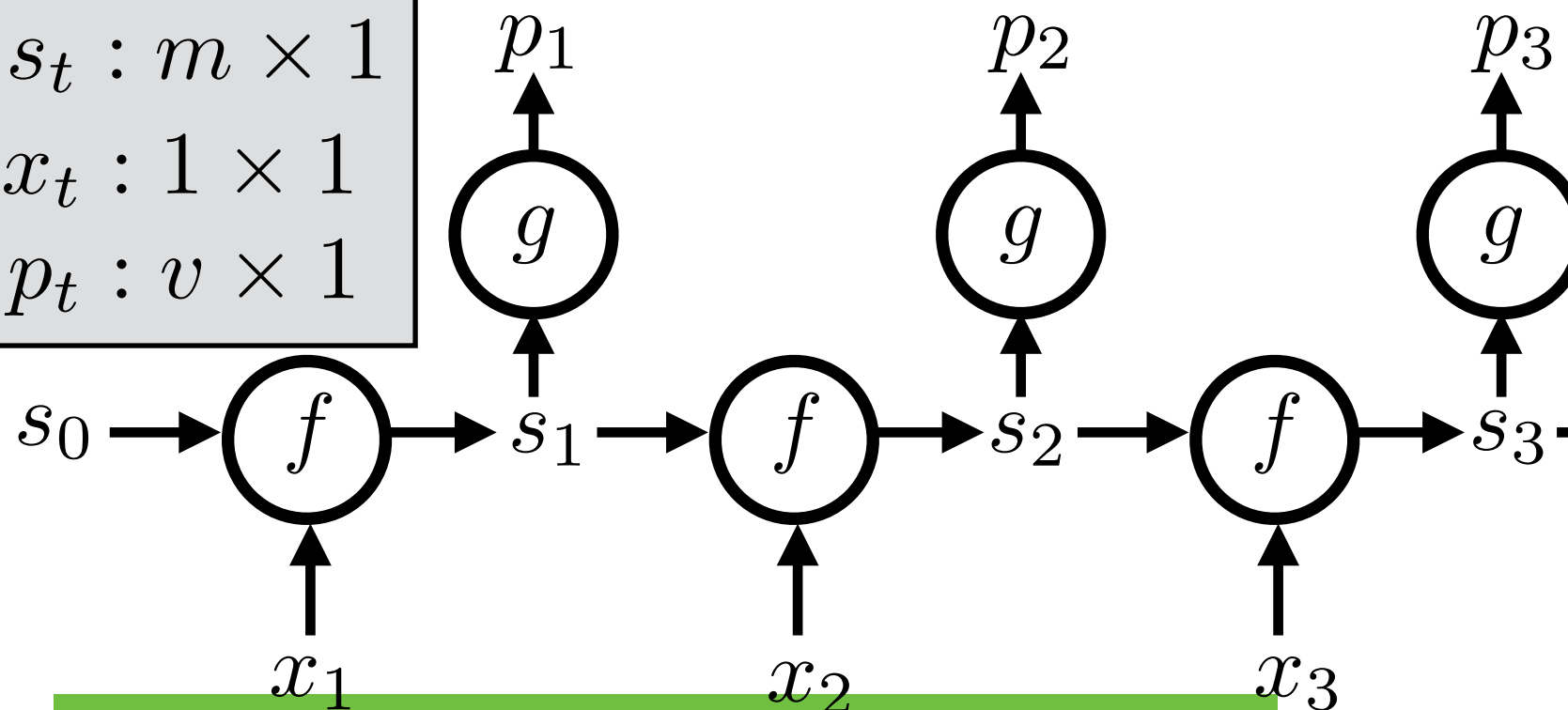
$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(W^o s_t + W_0^o) \end{aligned}$$

$v \times m \quad v \times 1$

$$\begin{aligned} L_{\text{seq}}(p^{(i)}, y^{(i)}) &= \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)}) \\ p^{(i)} &= R(x^{(i)}; W^o, W_0^o) \\ J(W^o, W_0^o) &= \sum_{i=1}^q L_{\text{seq}}(p^{(i)}, y^{(i)}) \end{aligned}$$

Can express as a state machine

$$\begin{aligned} s_t &: m \times 1 \\ x_t &: 1 \times 1 \\ p_t &: v \times 1 \end{aligned}$$



- Example: Alphabet $\{0,1\}$; state is last $m = 3$ characters

Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

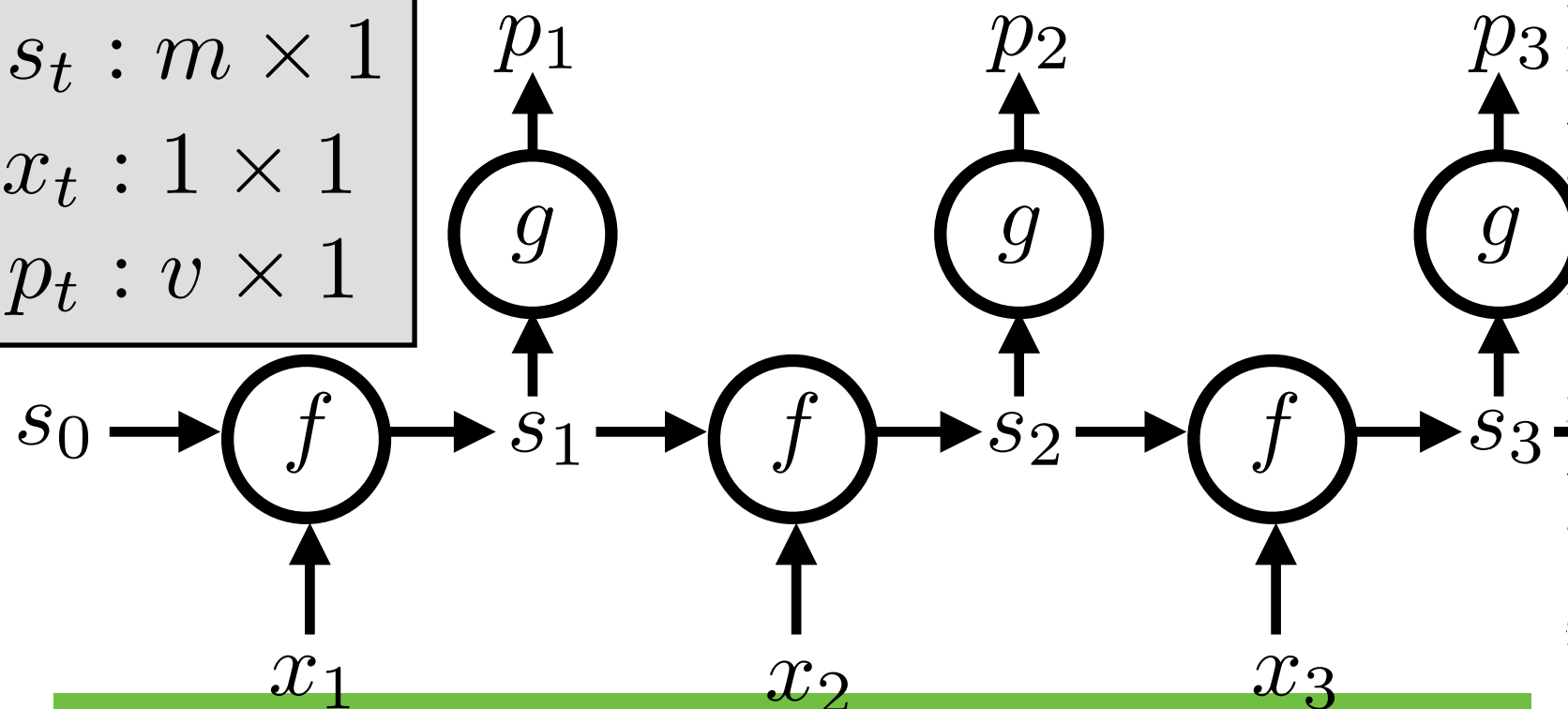
$$s_t = f_1(W^{sx}x_t + W^{ss}s_{t-1} + W_0^{ss})$$

$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(\underbrace{W^o}_{v \times m} s_t + \underbrace{W_0^o}_{v \times 1}) \end{aligned}$$

$$\begin{aligned} L_{\text{seq}}(p^{(i)}, y^{(i)}) &= \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)}) \\ p^{(i)} &= R(x^{(i)}; W^o, W_0^o) \\ J(W^o, W_0^o) &= \sum_{i=1}^q L_{\text{seq}}(p^{(i)}, y^{(i)}) \end{aligned}$$

Can express as a state machine

$$\begin{aligned} s_t &: m \times 1 \\ x_t &: 1 \times 1 \\ p_t &: v \times 1 \end{aligned}$$



- Example: Alphabet of ℓ characters; state is last c characters

Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

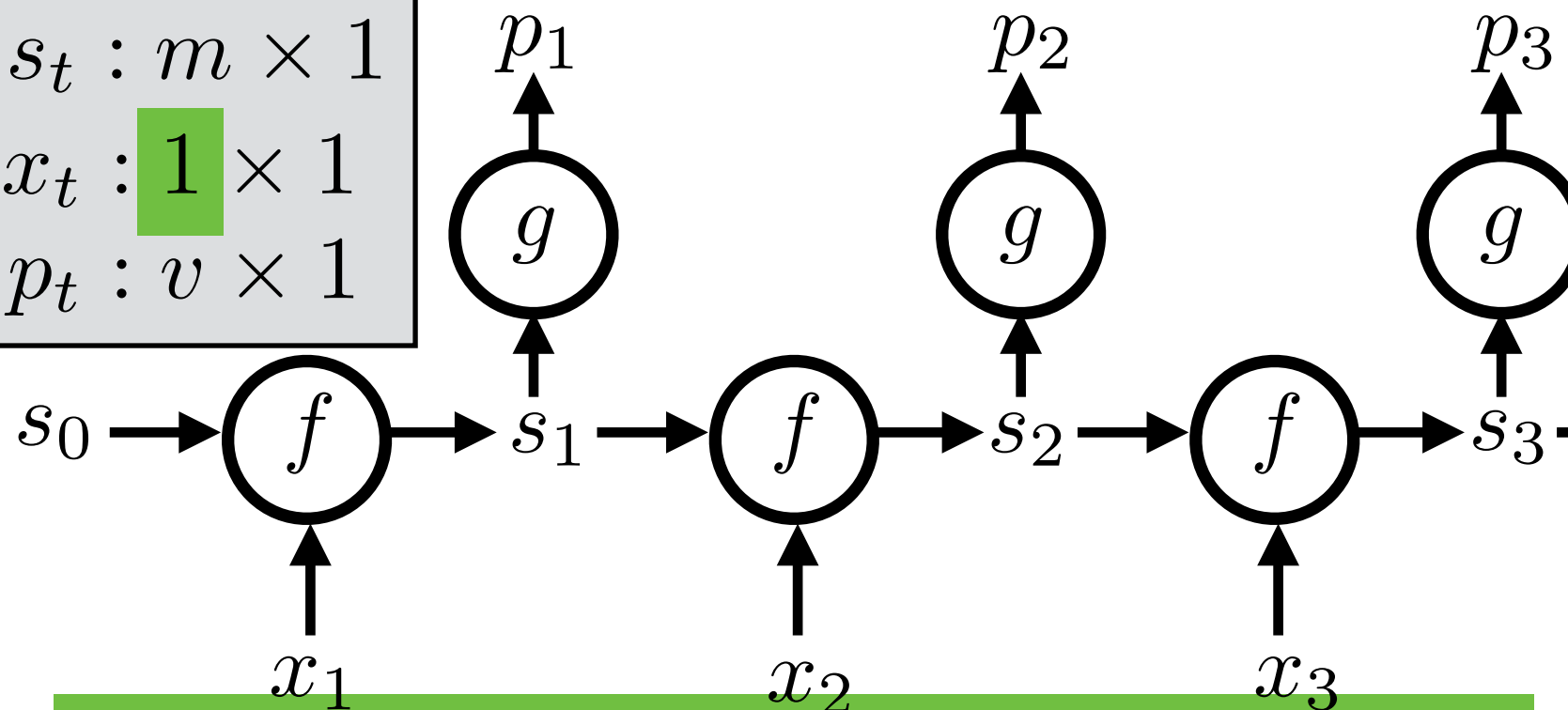
$$s_t = f_1(W^{sx}x_t + W^{ss}s_{t-1} + W_0^{ss})$$

$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(\underbrace{W^o}_{v \times m} s_t + \underbrace{W_0^o}_{v \times 1}) \end{aligned}$$

$$\begin{aligned} L_{\text{seq}}(p^{(i)}, y^{(i)}) &= \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)}) \\ p^{(i)} &= R(x^{(i)}; W^o, W_0^o) \\ J(W^o, W_0^o) &= \sum_{i=1}^q L_{\text{seq}}(p^{(i)}, y^{(i)}) \end{aligned}$$

Can express as a state machine

$$\begin{aligned} s_t &: m \times 1 \\ x_t &: 1 \times 1 \\ p_t &: v \times 1 \end{aligned}$$



- Example: Alphabet of ℓ characters; state is last c characters

Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

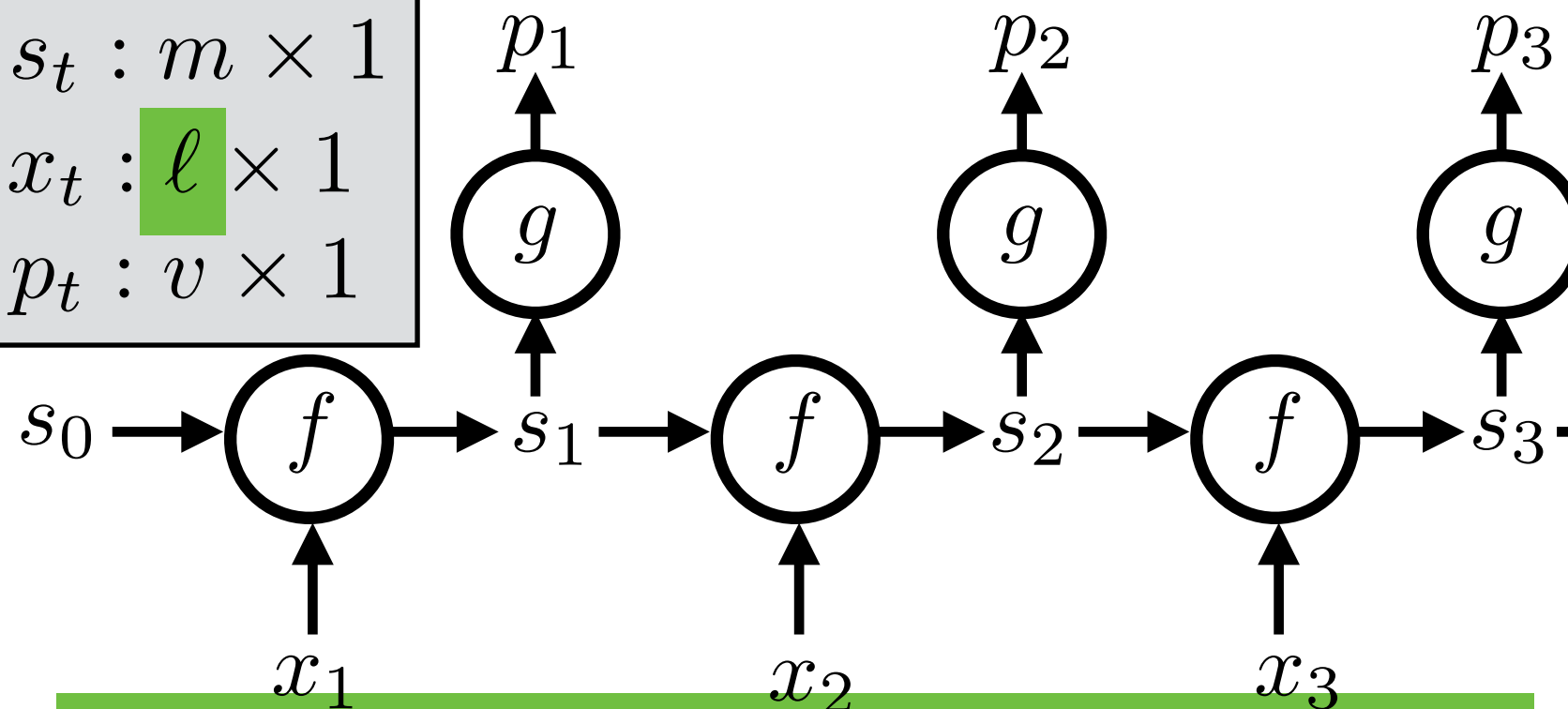
$$s_t = f_1(W^{sx}x_t + W^{ss}s_{t-1} + W_0^{ss})$$

$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(\underbrace{W^o}_{v \times m} s_t + \underbrace{W_0^o}_{v \times 1}) \end{aligned}$$

$$\begin{aligned} L_{\text{seq}}(p^{(i)}, y^{(i)}) &= \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)}) \\ p^{(i)} &= R(x^{(i)}; W^o, W_0^o) \\ J(W^o, W_0^o) &= \sum_{i=1}^q L_{\text{seq}}(p^{(i)}, y^{(i)}) \end{aligned}$$

Can express as a state machine

$$\begin{aligned} s_t &: m \times 1 \\ x_t &: \ell \times 1 \\ p_t &: v \times 1 \end{aligned}$$



- Example: Alphabet of ℓ characters; state is last c characters

Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

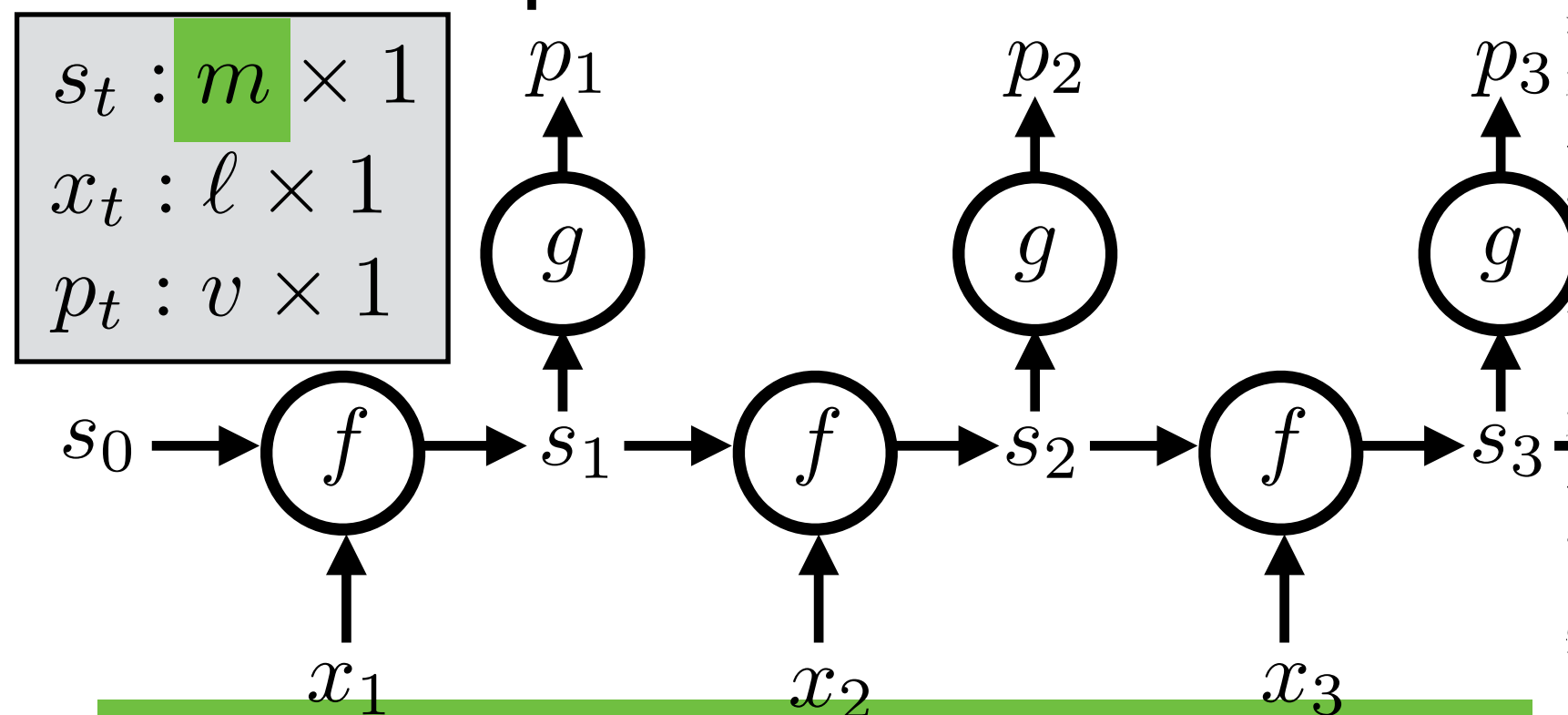
$$s_t = f_1 (W^{sx} x_t + W^{ss} s_{t-1} + W_0^{ss})$$

$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(W^o s_t + W_0^o) \end{aligned}$$

$v \times m \quad v \times 1$

$$\begin{aligned} L_{\text{seq}}(p^{(i)}, y^{(i)}) &= \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)}) \\ p^{(i)} &= R(x^{(i)}; W^o, W_0^o) \\ J(W^o, W_0^o) &= \sum_{i=1}^q L_{\text{seq}}(p^{(i)}, y^{(i)}) \end{aligned}$$

Can express as a state machine



- Example: Alphabet of ℓ characters; state is last c characters

Recall: familiar pattern

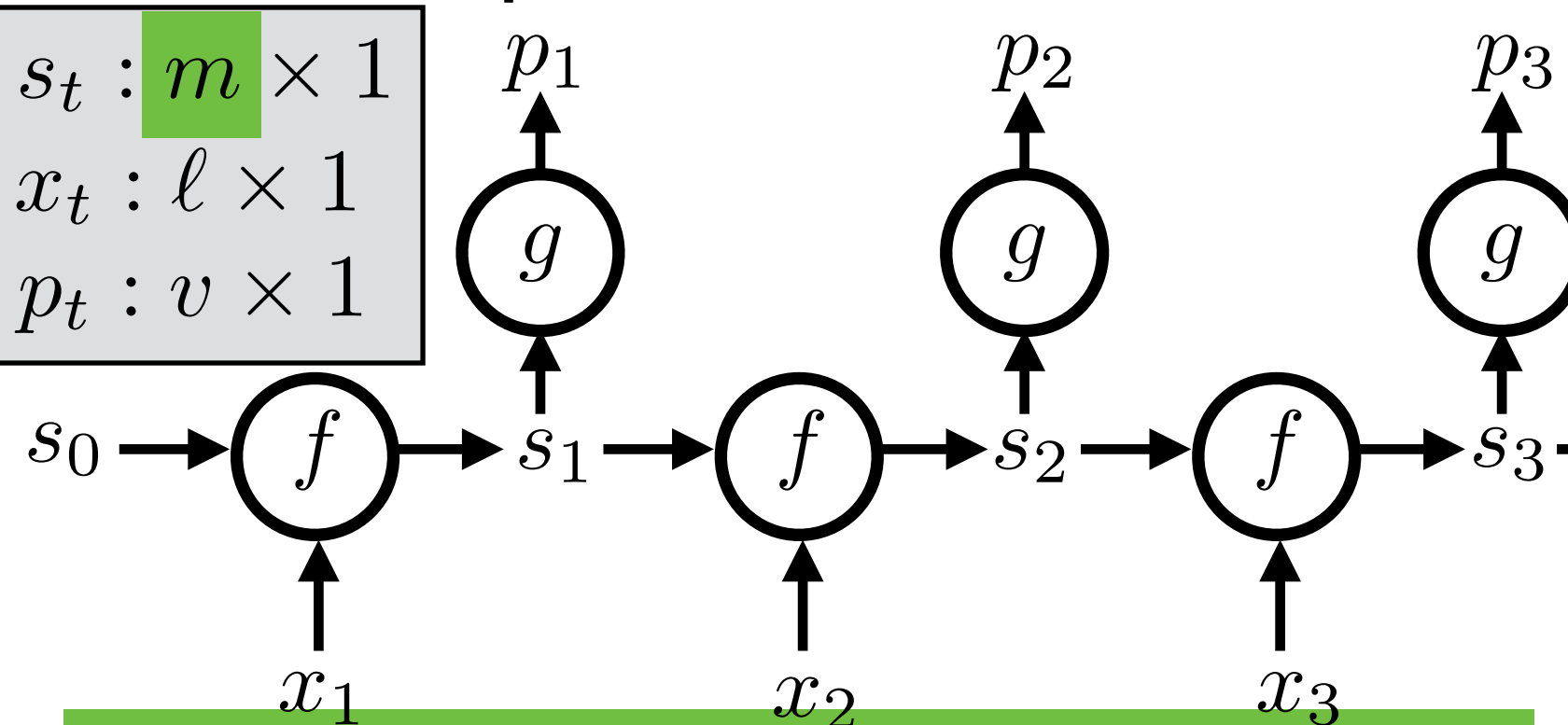
1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

$$s_t = f_1(W^{sx}x_t + W^{ss}s_{t-1} + W_0^{ss})$$

$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(\underbrace{W^o}_{v \times m} s_t + \underbrace{W_0^o}_{v \times 1}) \end{aligned}$$

$$\begin{aligned} L_{\text{seq}}(p^{(i)}, y^{(i)}) &= \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)}) \\ p^{(i)} &= R(x^{(i)}; W^o, W_0^o) \\ J(W^o, W_0^o) &= \sum_{i=1}^q L_{\text{seq}}(p^{(i)}, y^{(i)}) \end{aligned}$$

Can express as a state machine



- Example: Alphabet of ℓ characters; state is last c characters ($m = c\ell$)

Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

$$s_t = f_1(W^{sx}x_t + W^{ss}s_{t-1} + W_0^{ss})$$

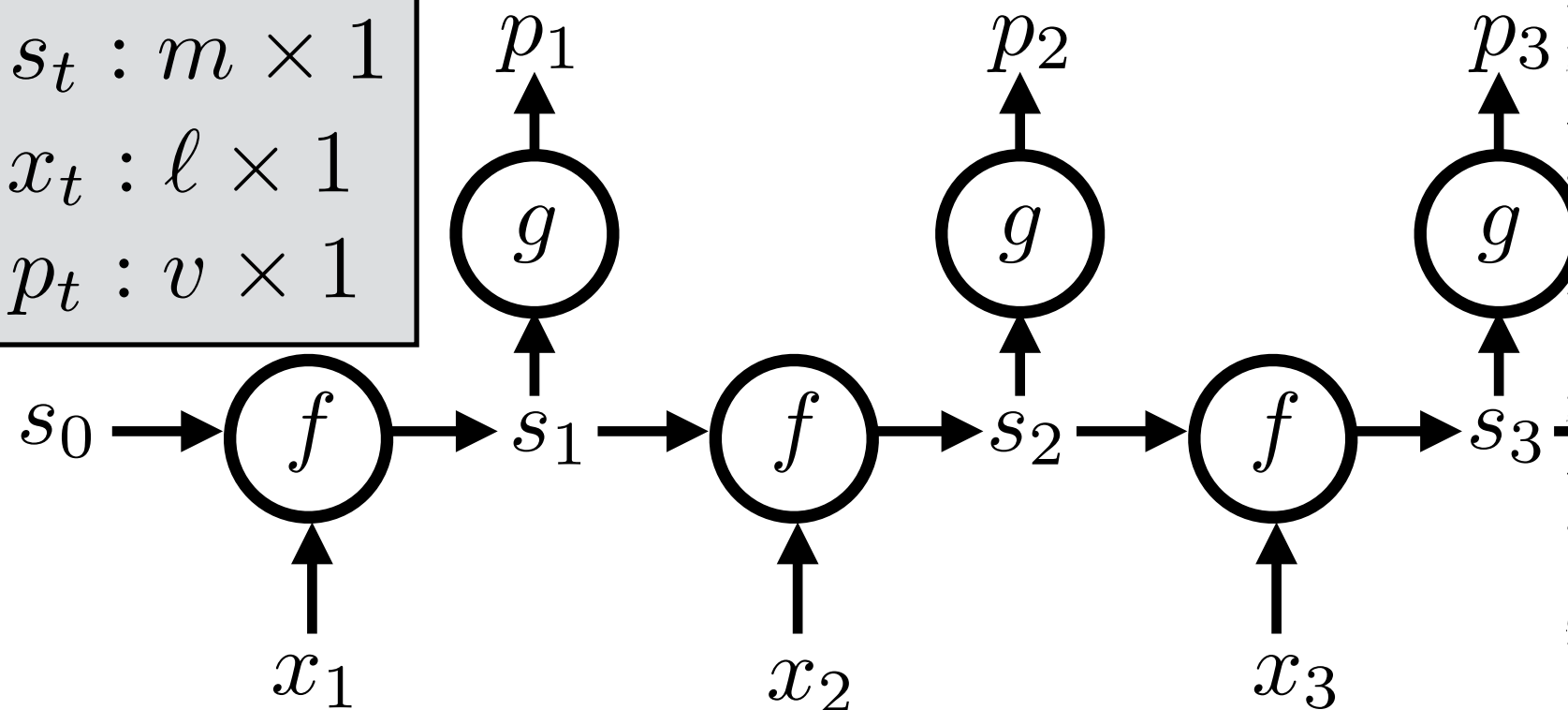
$$\begin{aligned}
 p_t &= g(s_t) \\
 &= f_2(W^o s_t + W_0^o)
 \end{aligned}$$

$v \times m \quad v \times 1$

$$\begin{aligned}
 L_{\text{seq}}(p^{(i)}, y^{(i)}) &= \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)}) \\
 p^{(i)} &= R(x^{(i)}; W^o, W_0^o) \\
 J(W^o, W_0^o) &= \sum_{i=1}^q L_{\text{seq}}(p^{(i)}, y^{(i)})
 \end{aligned}$$

Can express as a state machine

$$\begin{aligned} s_t &: m \times 1 \\ x_t &: \ell \times 1 \\ p_t &: v \times 1 \end{aligned}$$



- Example: Alphabet of ℓ characters; state is last c characters ($m = c\ell$)

Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

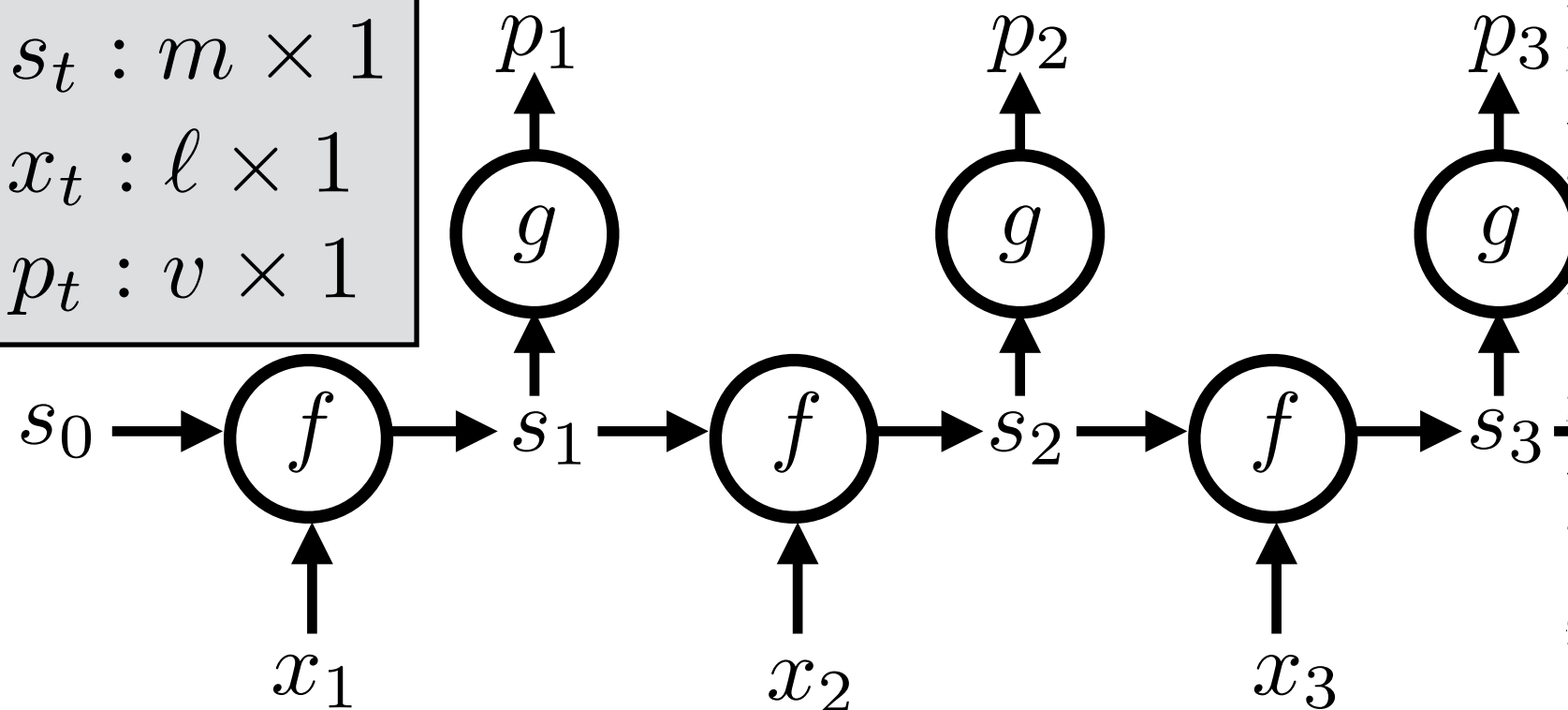
$$s_t = f_1(W^{sx}x_t + W^{ss}s_{t-1} + W_0^{ss})$$

$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(\underbrace{W^o}_{v \times m} s_t + \underbrace{W_0^o}_{v \times 1}) \end{aligned}$$

$$\begin{aligned} L_{\text{seq}}(p^{(i)}, y^{(i)}) &= \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)}) \\ p^{(i)} &= R(x^{(i)}; W^o, W_0^o) \\ J(W^o, W_0^o) &= \sum_{i=1}^q L_{\text{seq}}(p^{(i)}, y^{(i)}) \end{aligned}$$

Can express as a state machine

$$\begin{aligned} s_t &: m \times 1 \\ x_t &: \ell \times 1 \\ p_t &: v \times 1 \end{aligned}$$



- Example: Alphabet of ℓ characters; state is last c characters ($m = c\ell$)

- Recall: familiar pattern
1. Choose how to predict label (given features & parameters)
 2. Choose a loss (between guess & actual label)
 3. Choose parameters by trying to minimize the training loss

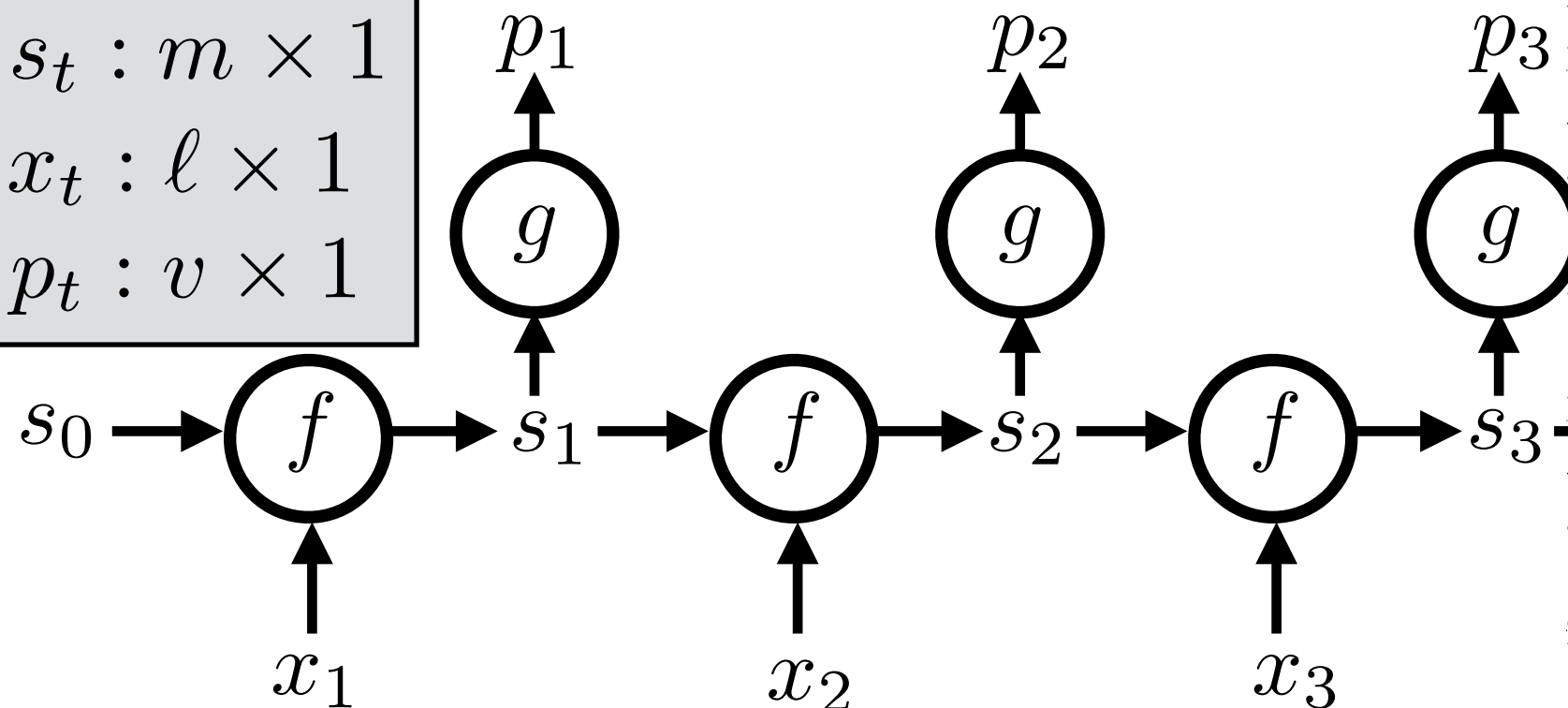
$$s_t = f_1(W^{sx}x_t + W^{ss}s_{t-1} + W_0^{ss})$$

$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(\underbrace{W^o}_{v \times m} s_t + \underbrace{W_0^o}_{v \times 1}) \end{aligned}$$

$$\begin{aligned} L_{\text{seq}}(p^{(i)}, y^{(i)}) &= \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)}) \\ p^{(i)} &= R(x^{(i)}; W^o, W_0^o) \\ J(W^o, W_0^o) &= \sum_{i=1}^q L_{\text{seq}}(p^{(i)}, y^{(i)}) \end{aligned}$$

Can express as a state machine

$$\begin{aligned} s_t &: m \times 1 \\ x_t &: \ell \times 1 \\ p_t &: v \times 1 \end{aligned}$$



- Example: Alphabet of ℓ characters; state is last c characters ($m = c\ell$)

Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

$$s_t = f_1 (W^{sx} x_t + W^{ss} s_{t-1} + W_0^{ss})$$

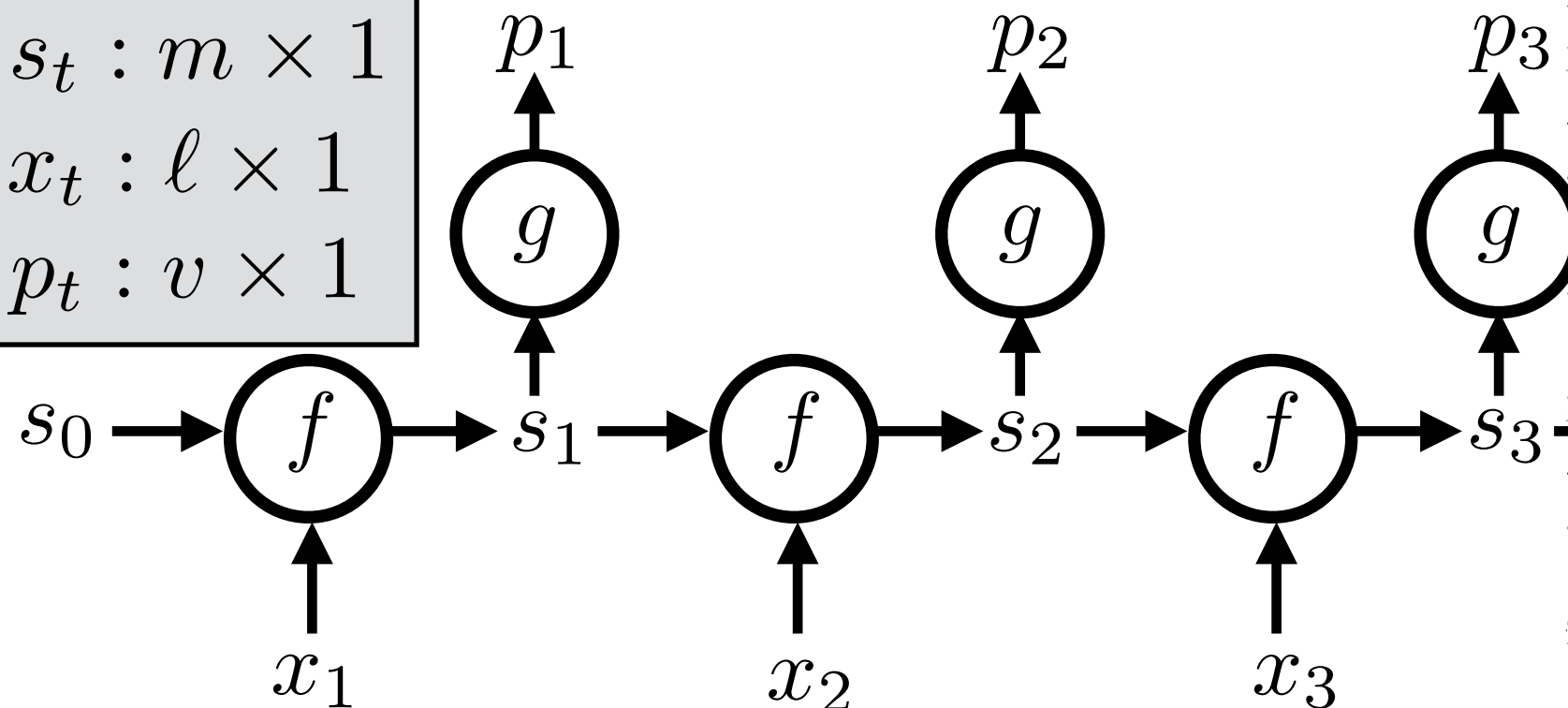
$$\begin{aligned} p_t &= g(s_t) \\ &= f_2 (W^o s_t + W_0^o) \end{aligned}$$

$v \times m \quad v \times 1$

$$\begin{aligned} L_{\text{seq}}(p^{(i)}, y^{(i)}) &= \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)}) \\ p^{(i)} &= R(x^{(i)}; W^o, W_0^o) \\ J(W^o, W_0^o) &= \sum_{i=1}^q L_{\text{seq}}(p^{(i)}, y^{(i)}) \end{aligned}$$

Can express as a state machine

$$\begin{aligned} s_t &: m \times 1 \\ x_t &: \ell \times 1 \\ p_t &: v \times 1 \end{aligned}$$



- Example: Alphabet of ℓ characters; state is last c characters ($m = c\ell$)

Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

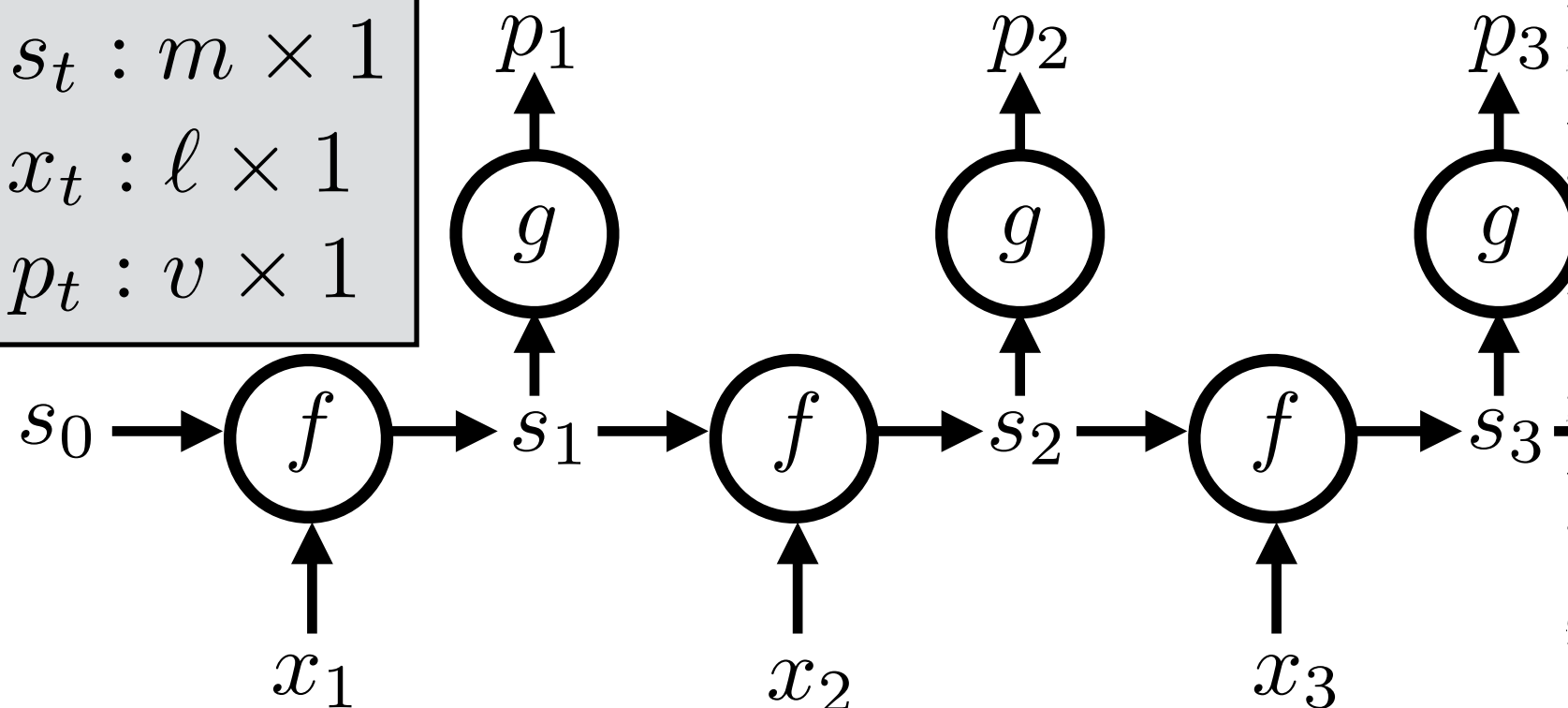
$$s_t = f_1(W^{sx}x_t + W^{ss}s_{t-1} + W_0^{ss})$$

$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(\underbrace{W^o}_{v \times m} s_t + \underbrace{W_0^o}_{v \times 1}) \end{aligned}$$

$$\begin{aligned} L_{\text{seq}}(p^{(i)}, y^{(i)}) &= \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)}) \\ p^{(i)} &= R(x^{(i)}; W^o, W_0^o) \\ J(W^o, W_0^o) &= \sum_{i=1}^q L_{\text{seq}}(p^{(i)}, y^{(i)}) \end{aligned}$$

Can express as a state machine

$$\begin{aligned} s_t &: m \times 1 \\ x_t &: \ell \times 1 \\ p_t &: v \times 1 \end{aligned}$$



- Example: Alphabet of ℓ characters; state is last c characters ($m = c\ell$)

Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

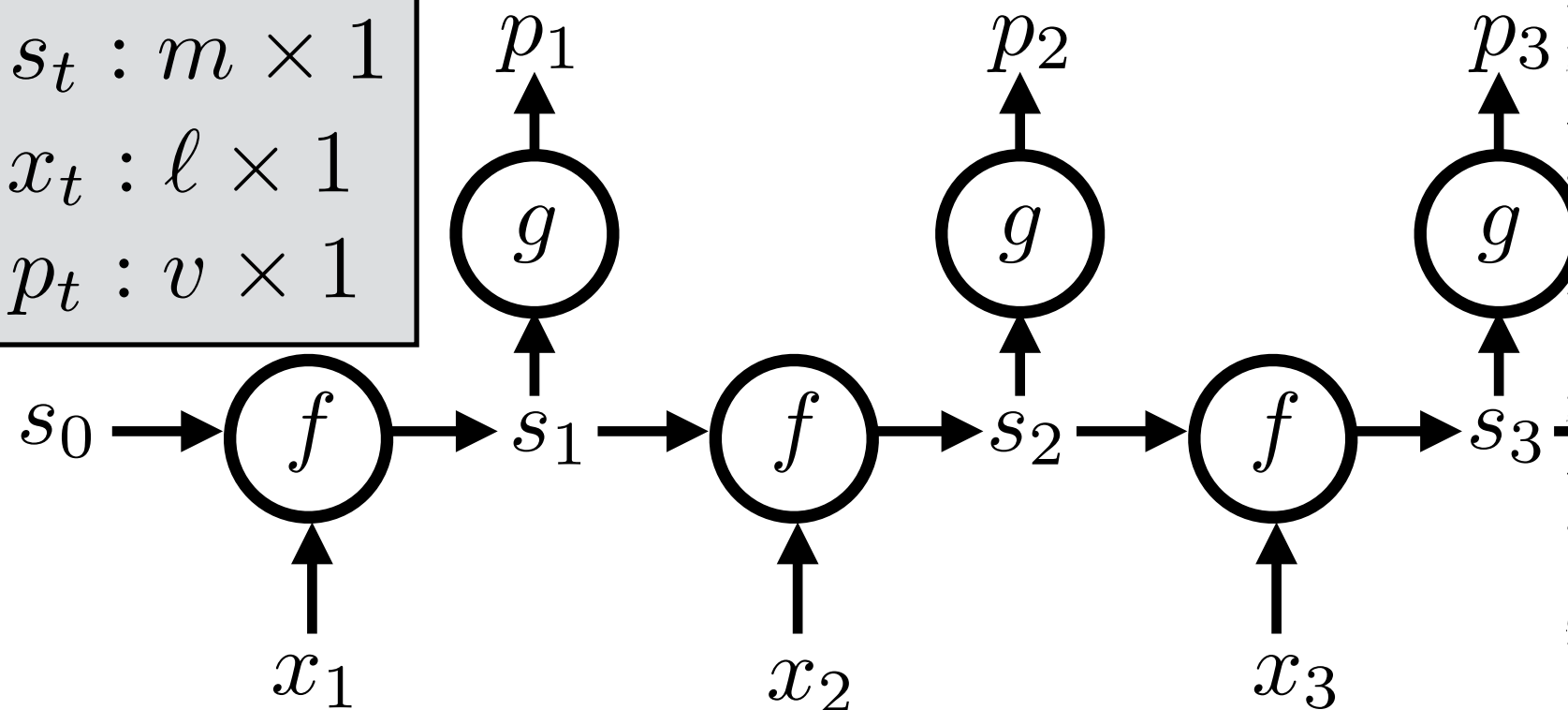
$$s_t = f_1(W^{sx}x_t + W^{ss}s_{t-1} + W_0^{ss})$$

$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(\underbrace{W^o}_{v \times m} s_t + \underbrace{W_0^o}_{v \times 1}) \end{aligned}$$

$$\begin{aligned} L_{\text{seq}}(p^{(i)}, y^{(i)}) &= \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)}) \\ p^{(i)} &= \text{R}(x^{(i)}; W, W_0) \\ J(W, W_0) &= \sum_{i=1}^q L_{\text{seq}}(p^{(i)}, y^{(i)}) \end{aligned}$$

Can express as a state machine

$$\begin{aligned} s_t &: m \times 1 \\ x_t &: \ell \times 1 \\ p_t &: v \times 1 \end{aligned}$$



- Example: Alphabet of ℓ characters; state is last c characters ($m = c\ell$)

Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

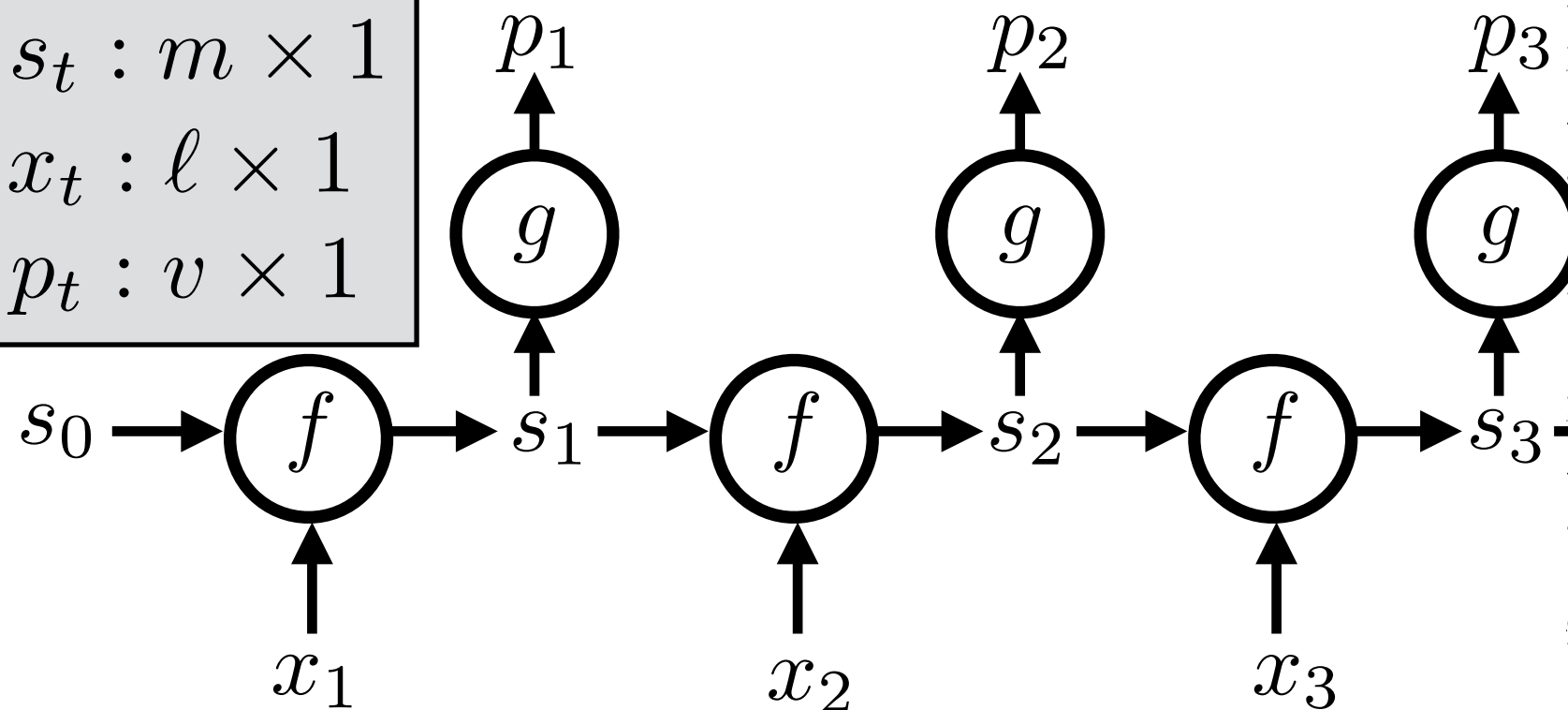
$$s_t = f_1(W^{sx}x_t + W^{ss}s_{t-1} + W_0^{ss})$$

$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(\underbrace{W^o}_{v \times m} s_t + \underbrace{W_0^o}_{v \times 1}) \end{aligned}$$

$$\begin{aligned} L_{\text{seq}}(p^{(i)}, y^{(i)}) &= \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)}) \\ p^{(i)} &= \mathbf{R}(x^{(i)}; W, W_0) \\ J(W, W_0) &= \sum_{i=1}^q L_{\text{seq}}(p^{(i)}, y^{(i)}) \end{aligned}$$

Can express as a state machine

$$\begin{aligned} s_t &: m \times 1 \\ x_t &: \ell \times 1 \\ p_t &: v \times 1 \end{aligned}$$



- Example: Alphabet of ℓ characters; state is last c characters ($m = c\ell$)

Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

$$s_t = f_1(W^{sx}x_t + W^{ss}s_{t-1} + W_0^{ss})$$

$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(\underbrace{W^o}_{v \times m} s_t + \underbrace{W_0^o}_{v \times 1}) \end{aligned}$$

$$\begin{aligned} L_{\text{seq}}(p^{(i)}, y^{(i)}) &= \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)}) \\ p^{(i)} &= \text{RNN}(x^{(i)}; W, W_0) \\ J(W, W_0) &= \sum_{i=1}^q L_{\text{seq}}(p^{(i)}, y^{(i)}) \end{aligned}$$

Recurrent neural network

Recurrent neural network

- Example: Alphabet of ℓ chars; state is last c chars ($m = c\ell$)

Recurrent neural network

- Example: Alphabet of ℓ chars; state is last c chars ($m = c\ell$)

$$s_t = f_1 (W^{sx} x_t + W^{ss} s_{t-1} + W_0^{ss})$$

$$p_t = f_2(W^o s_t + W_0^o)$$

Recurrent neural network

- Example: Alphabet of ℓ chars; state is last c chars ($m = c\ell$)
 $s_t = f_1 (W^{sx} x_t + W^{ss} s_{t-1} + W_0^{ss})$
 $p_t = f_2(W^o s_t + W_0^o)$ $p^{(i)} = \text{RNN}(x^{(i)}; W, W_0)$

Recurrent neural network

$$s_t : m \times 1$$

$$x_t : \ell \times 1$$

$$p_t : v \times 1$$

- Example: Alphabet of ℓ chars; state is last c chars ($m = c\ell$)

$$s_t = f_1 (W^{sx} x_t + W^{ss} s_{t-1} + W_0^{ss})$$

$$p_t = f_2 (W^o s_t + W_0^o)$$

$$p^{(i)} = \text{RNN}(x^{(i)}; W, W_0)$$

Recurrent neural network

$$s_t : m \times 1$$

$$x_t : \ell \times 1$$

$$p_t : v \times 1$$

$$s_0$$

$$x_1$$

- Example: Alphabet of ℓ chars; state is last c chars ($m = c\ell$)

$$s_t = f_1(W^{sx}x_t + W^{ss}s_{t-1} + W_0^{ss})$$

$$p_t = f_2(W^o s_t + W_0^o)$$

$$p^{(i)} = \text{RNN}(x^{(i)}; W, W_0)$$

Recurrent neural network

$$s_t : m \times 1$$

$$x_t : \ell \times 1$$

$$p_t : v \times 1$$

$s_0 \rightarrow$

\uparrow
 x_1

- Example: Alphabet of ℓ chars; state is last c chars ($m = c\ell$)

$$s_t = f_1(W^{sx}x_t + W^{ss}s_{t-1} + W_0^{ss})$$

$$p_t = f_2(W^o s_t + W_0^o)$$

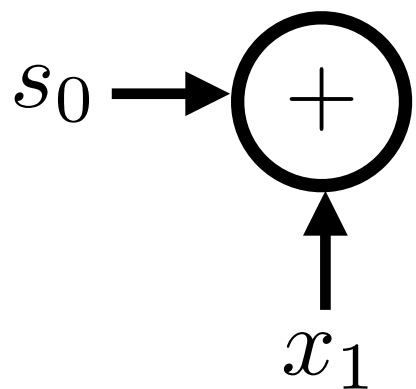
$$p^{(i)} = \text{RNN}(x^{(i)}; W, W_0)$$

Recurrent neural network

$$s_t : m \times 1$$

$$x_t : \ell \times 1$$

$$p_t : v \times 1$$



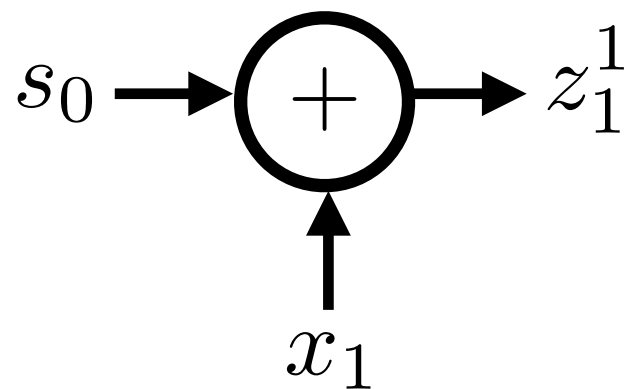
- Example: Alphabet of ℓ chars; state is last c chars ($m = c\ell$)
$$s_t = f_1(W^{sx}x_t + W^{ss}s_{t-1} + W_0^{ss})$$
$$p_t = f_2(W^os_t + W_0^o)$$
$$p^{(i)} = \text{RNN}(x^{(i)}; W, W_0)$$

Recurrent neural network

$$s_t : m \times 1$$

$$x_t : \ell \times 1$$

$$p_t : v \times 1$$



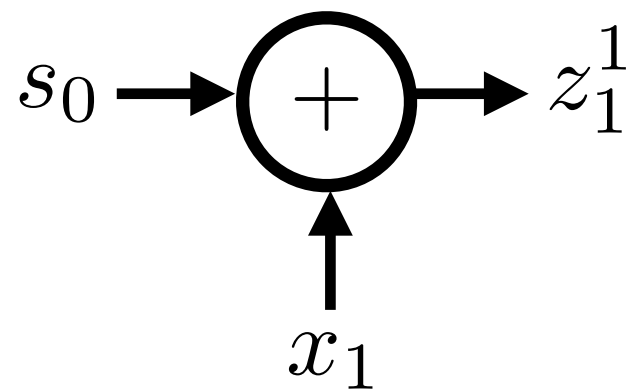
- Example: Alphabet of ℓ chars; state is last c chars ($m = c\ell$)
$$s_t = f_1 (W^{sx} x_t + W^{ss} s_{t-1} + W_0^{ss})$$
$$p_t = f_2 (W^o s_t + W_0^o)$$
$$p^{(i)} = \text{RNN}(x^{(i)}; W, W_0)$$

Recurrent neural network

$$s_t : m \times 1$$

$$x_t : \ell \times 1$$

$$p_t : v \times 1$$



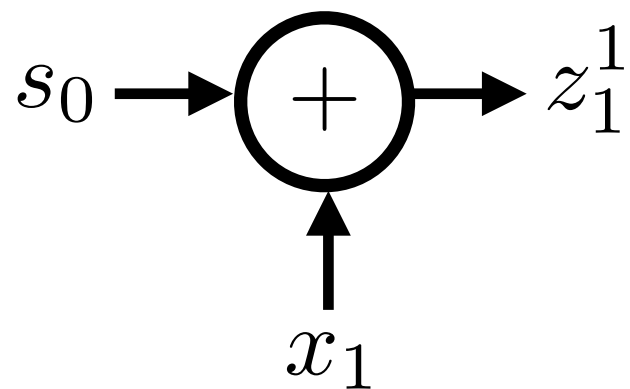
- Example: Alphabet of ℓ chars; state is last c chars ($m = c\ell$)
$$s_t = f_1(W^{sx}x_t + W^{ss}s_{t-1} + W_0^{ss})$$
$$p_t = f_2(W^os_t + W_0^o)$$
$$p^{(i)} = \text{RNN}(x^{(i)}; W, W_0)$$

Recurrent neural network

$$s_t : m \times 1$$

$$x_t : \ell \times 1$$

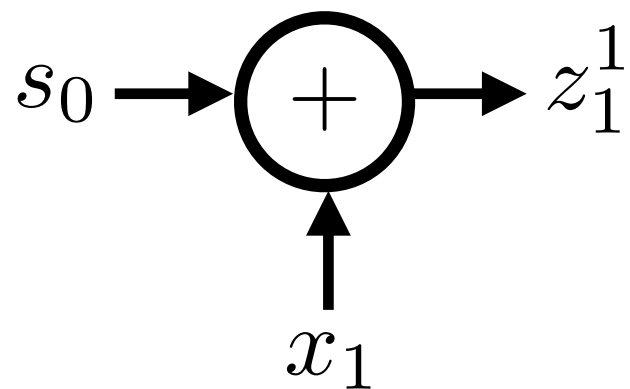
$$p_t : v \times 1$$



- Example: Alphabet of ℓ chars; state is last c chars ($m = c\ell$)
$$s_t = f_1(W^{sx}x_t + W^{ss}s_{t-1} + W_0^{ss})$$
$$p_t = f_2(W^os_t + W_0^o)$$
$$p^{(i)} = \text{RNN}(x^{(i)}; W, W_0)$$

Recurrent neural network

$$\begin{array}{l} s_t : m \times 1 \\ x_t : \ell \times 1 \\ p_t : v \times 1 \end{array}$$



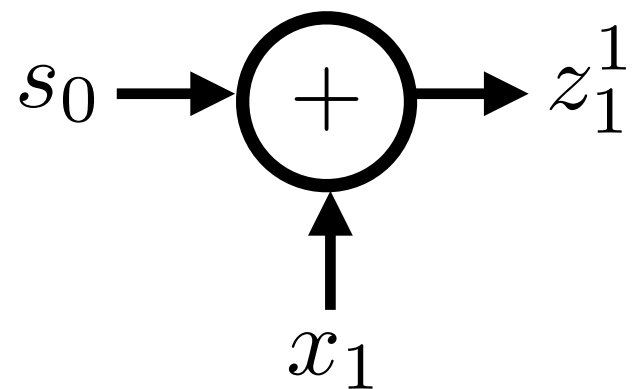
- Example: Alphabet of ℓ chars; state is last c chars ($m = c\ell$)
- $$s_t = f_1 (W^{sx} x_t + W^{ss} s_{t-1} + W_0^{ss})$$
- $$p_t = f_2 (W^o s_t + W_0^o) \rightarrow z_t^1$$
- $$p^{(i)} = \text{RNN}(x^{(i)}; W, W_0)$$

Recurrent neural network

$$s_t : m \times 1$$

$$x_t : \ell \times 1$$

$$p_t : v \times 1$$



- Example: Alphabet of ℓ chars; state is last c chars ($m = c\ell$)

$$s_t = f_1 \left(\underbrace{W^{sx} x_t + W^{ss} s_{t-1}}_{z_t^1} + W_0^{ss} \right)$$

$$p_t = f_2 \left(\underbrace{W^o s_t + W_0^o}_{z_t^2} \right)$$

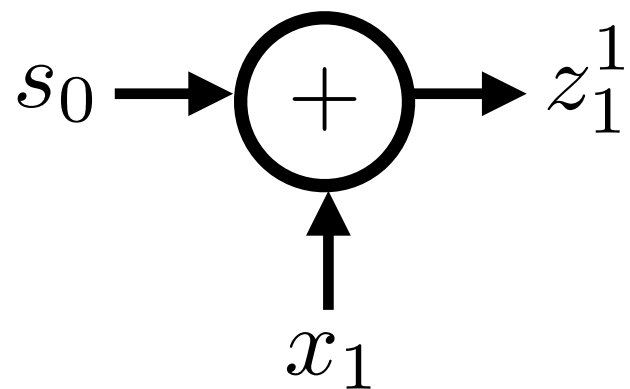
$$p^{(i)} = \text{RNN}(x^{(i)}; W, W_0)$$

Recurrent neural network

$$s_t : m \times 1$$

$$x_t : \ell \times 1$$

$$p_t : v \times 1$$



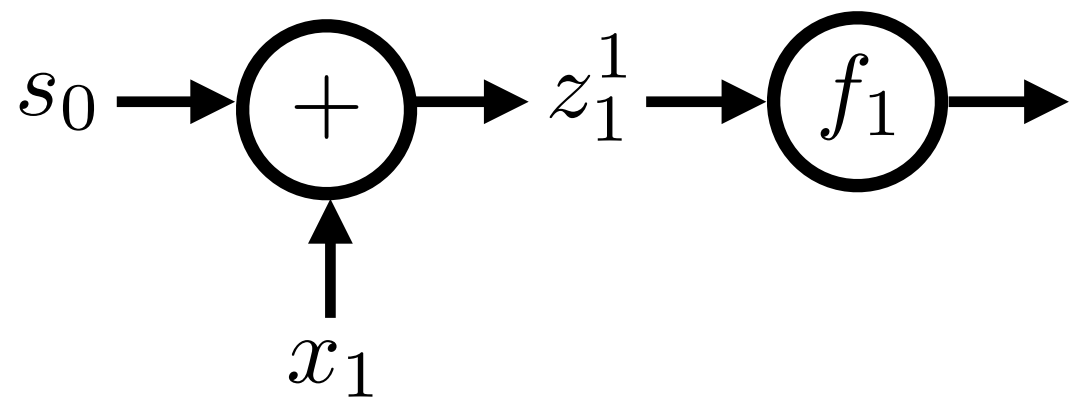
- Example: Alphabet of ℓ chars; state is last c chars ($m = c\ell$)
- $$s_t = f_1 \left(\underbrace{W^{sx} x_t + W^{ss} s_{t-1} + W_0^{ss}}_{z_t^1} \right)$$
- $$p_t = f_2 \left(\underbrace{W^o s_t + W_0^o}_{z_t^2} \right)$$
- $$p^{(i)} = \text{RNN}(x^{(i)}; W, W_0)$$

Recurrent neural network

$$s_t : m \times 1$$

$$x_t : \ell \times 1$$

$$p_t : v \times 1$$



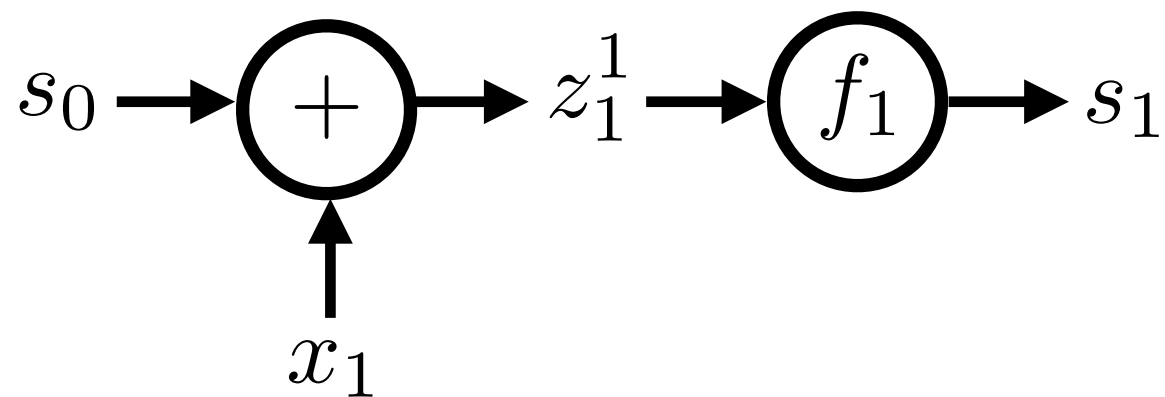
- Example: Alphabet of ℓ chars; state is last c chars ($m = c\ell$)
- $$s_t = f_1 \left(\underbrace{W^{sx} x_t + W^{ss} s_{t-1} + W_0^{ss}}_{z_t^1} \right)$$
- $$p_t = f_2 \left(\underbrace{W^o s_t + W_0^o}_{z_t^2} \right)$$
- $$p^{(i)} = \text{RNN}(x^{(i)}; W, W_0)$$

Recurrent neural network

$$s_t : m \times 1$$

$$x_t : \ell \times 1$$

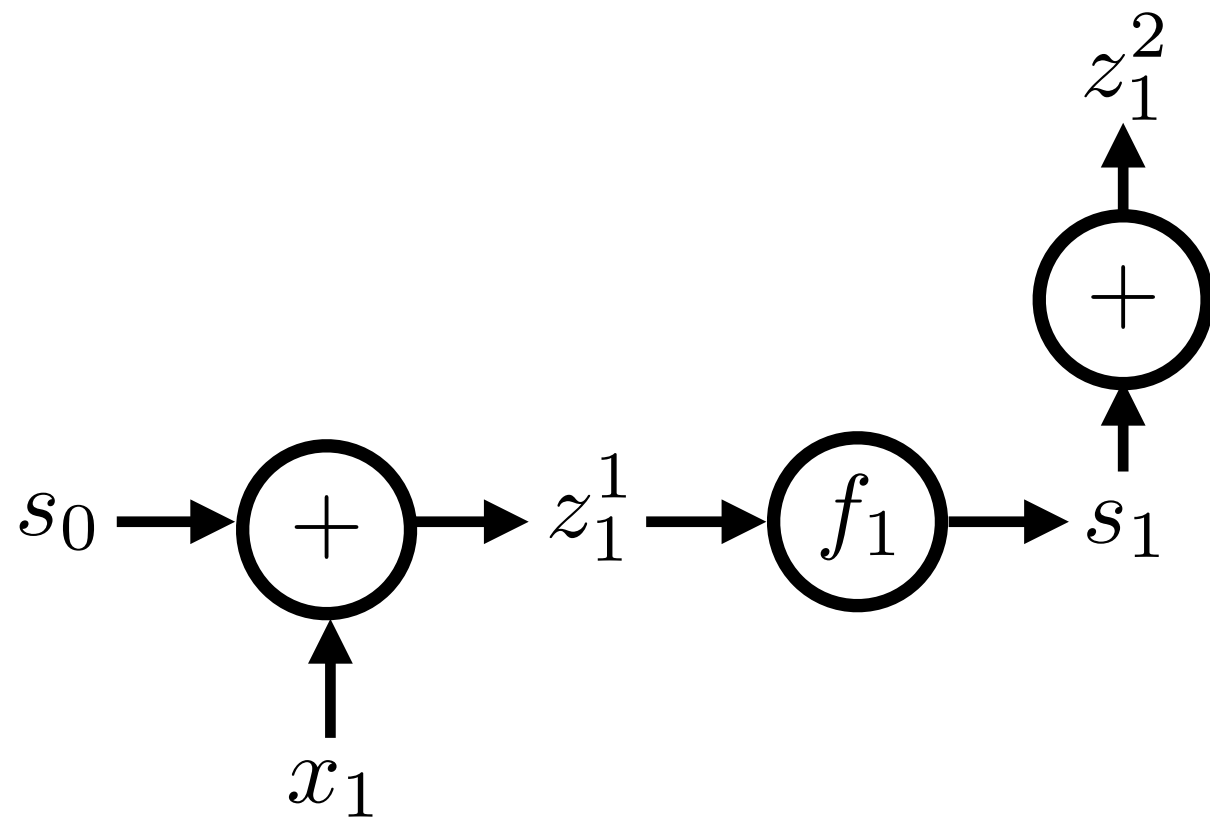
$$p_t : v \times 1$$



- Example: Alphabet of ℓ chars; state is last c chars ($m = c\ell$)
- $$s_t = f_1 \left(\underbrace{W^{sx} x_t + W^{ss} s_{t-1} + W_0^{ss}}_{z_t^1} \right)$$
- $$p_t = f_2 \left(\underbrace{W^o s_t + W_0^o}_{z_t^2} \right)$$
- $$p^{(i)} = \text{RNN}(x^{(i)}; W, W_0)$$

Recurrent neural network

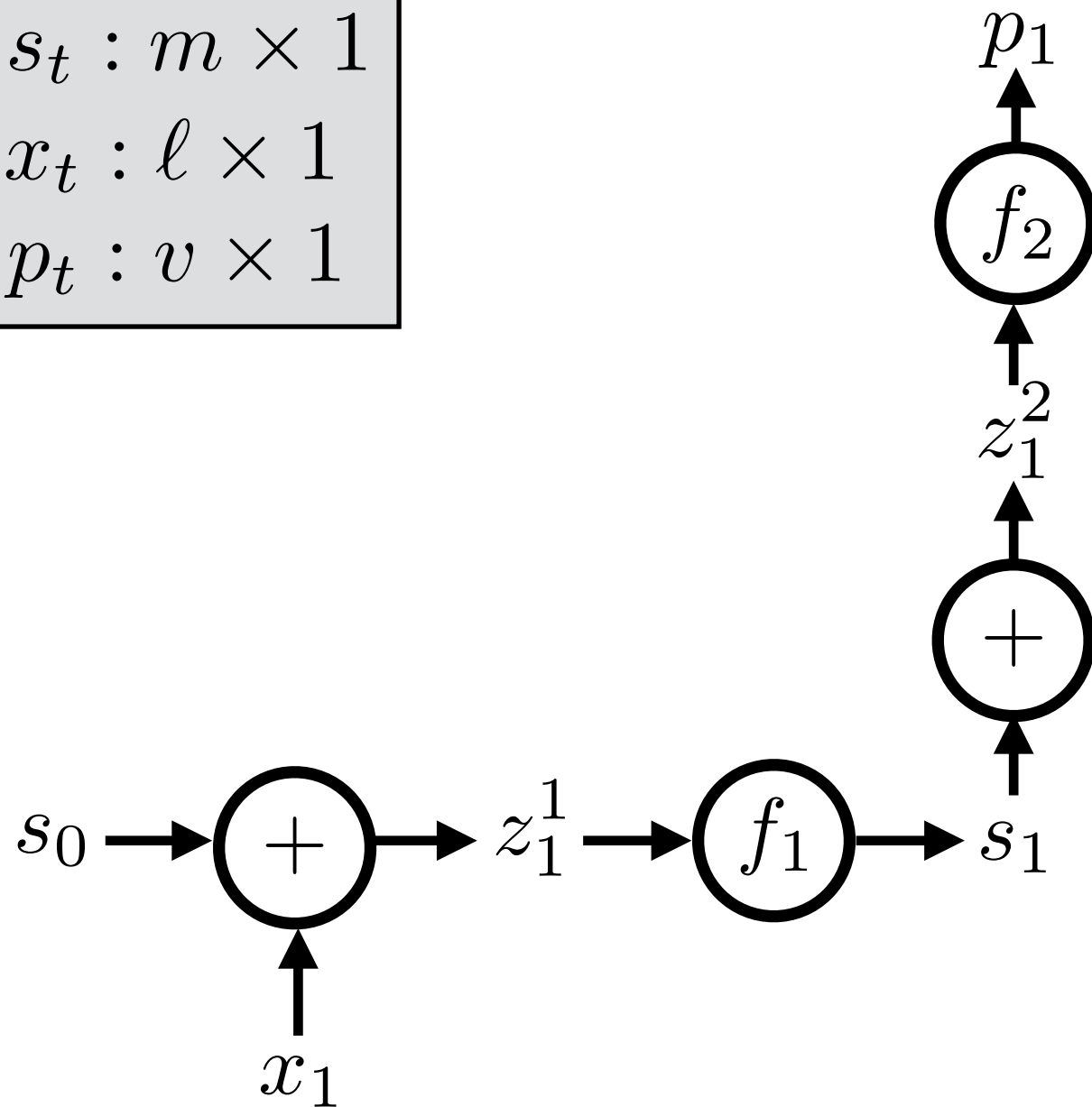
$$\begin{array}{l} s_t : m \times 1 \\ x_t : \ell \times 1 \\ p_t : v \times 1 \end{array}$$



- Example: Alphabet of ℓ chars; state is last c chars ($m = c\ell$)
- $$s_t = f_1 \left(\underbrace{W^{sx} x_t + W^{ss} s_{t-1} + W_0^{ss}}_{z_t^1} \right)$$
- $$p_t = f_2 \left(\underbrace{W^o s_t + W_0^o}_{z_t^2} \right)$$
- $$p^{(i)} = \text{RNN}(x^{(i)}; W, W_0)$$

Recurrent neural network

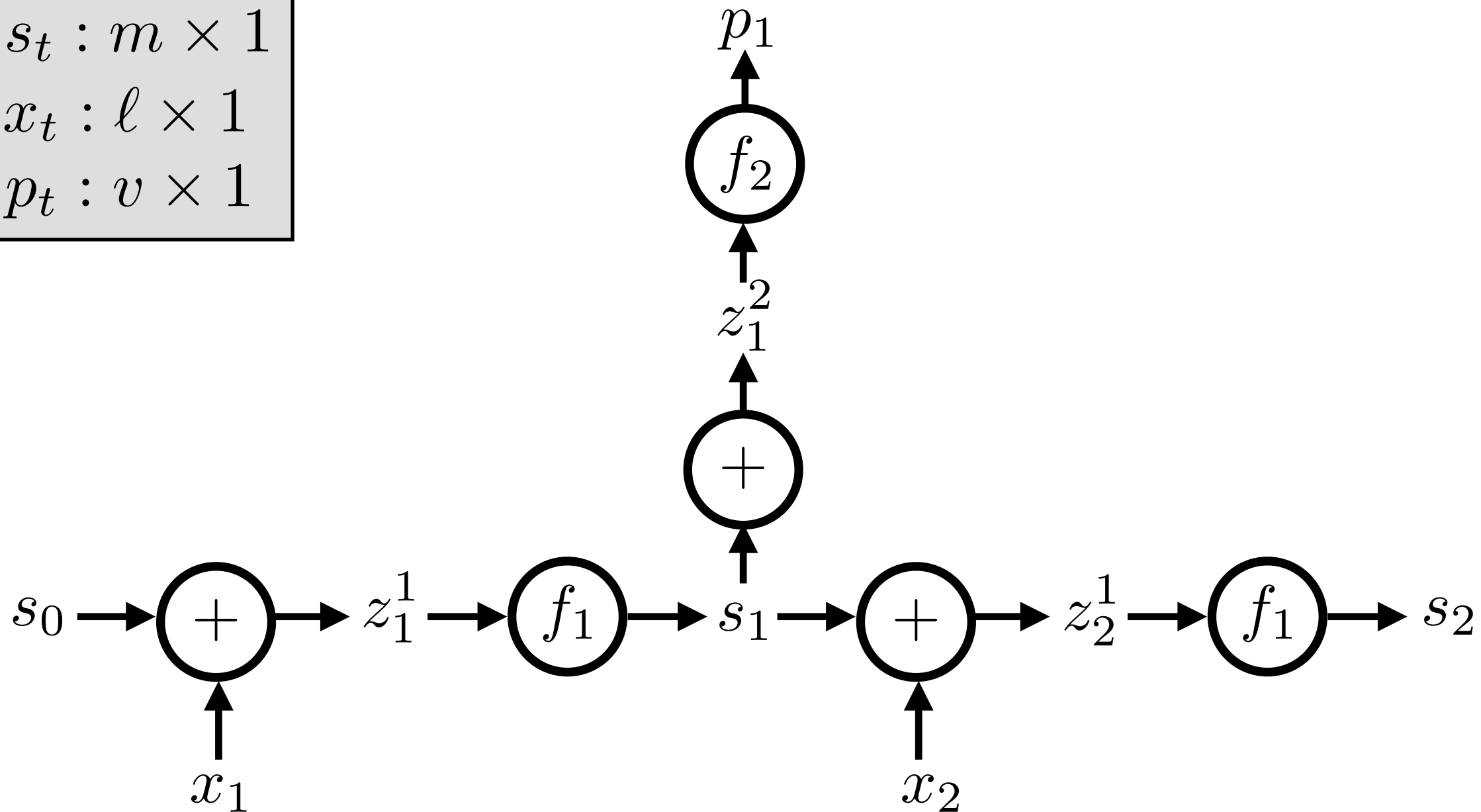
$$\begin{array}{l} s_t : m \times 1 \\ x_t : \ell \times 1 \\ p_t : v \times 1 \end{array}$$



- Example: Alphabet of ℓ chars; state is last c chars ($m = c\ell$)
- $$s_t = f_1 \left(\underbrace{W^{sx} x_t + W^{ss} s_{t-1} + W_0^{ss}}_{z_t^1} \right)$$
- $$p_t = f_2 \left(\underbrace{W^o s_t + W_0^o}_{z_t^2} \right)$$
- $$p^{(i)} = \text{RNN}(x^{(i)}; W, W_0)$$

Recurrent neural network

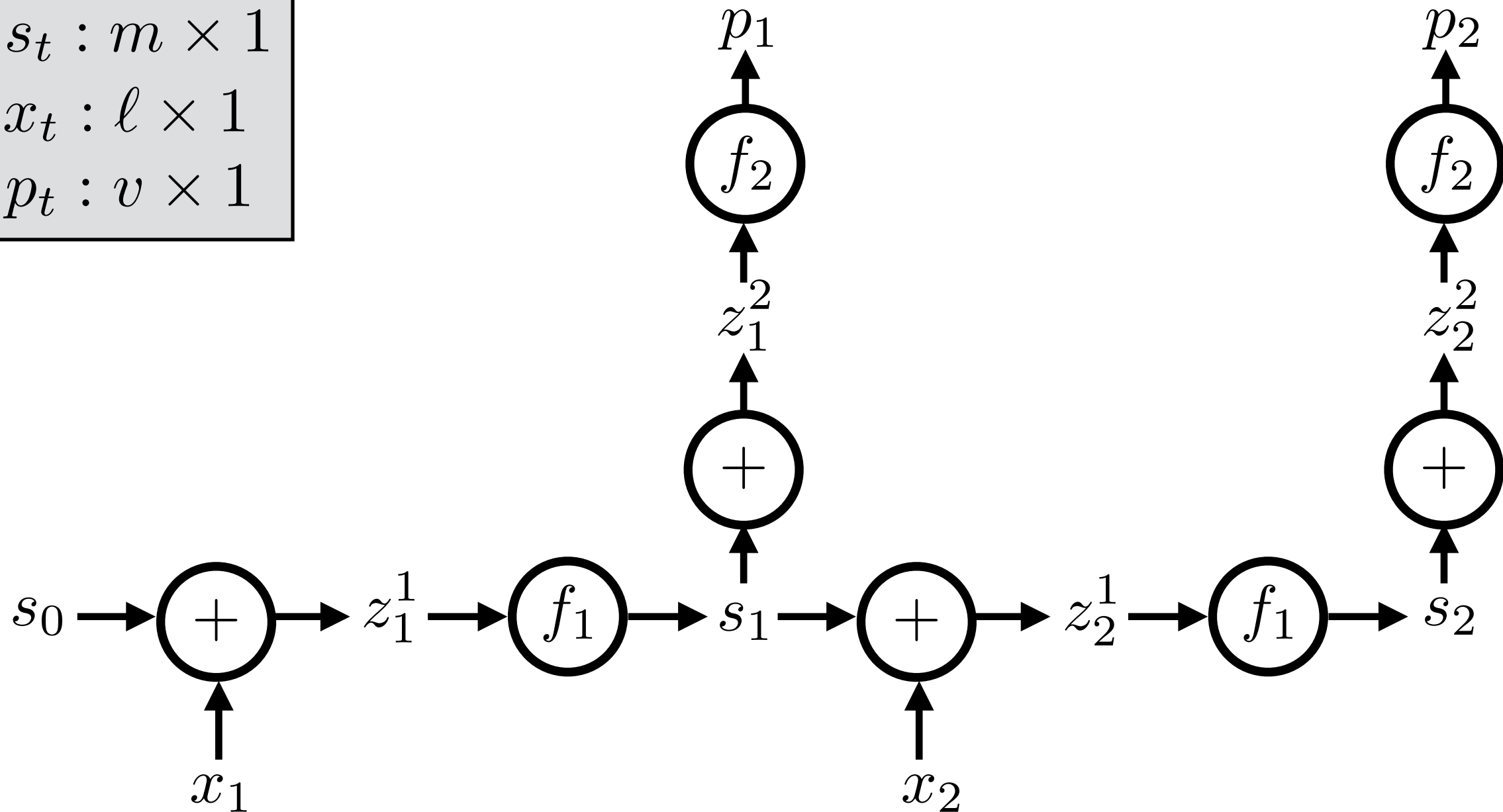
$$\begin{array}{l} s_t : m \times 1 \\ x_t : \ell \times 1 \\ p_t : v \times 1 \end{array}$$



- Example: Alphabet of ℓ chars; state is last c chars ($m = c\ell$)
- $$s_t = f_1 \left(\underbrace{W^{sx} x_t + W^{ss} s_{t-1} + W_0^{ss}}_{z_t^1} \right)$$
- $$p_t = f_2 \left(\underbrace{W^o s_t + W_0^o}_{z_t^2} \right)$$
- $$p^{(i)} = \text{RNN}(x^{(i)}; W, W_0)$$

Recurrent neural network

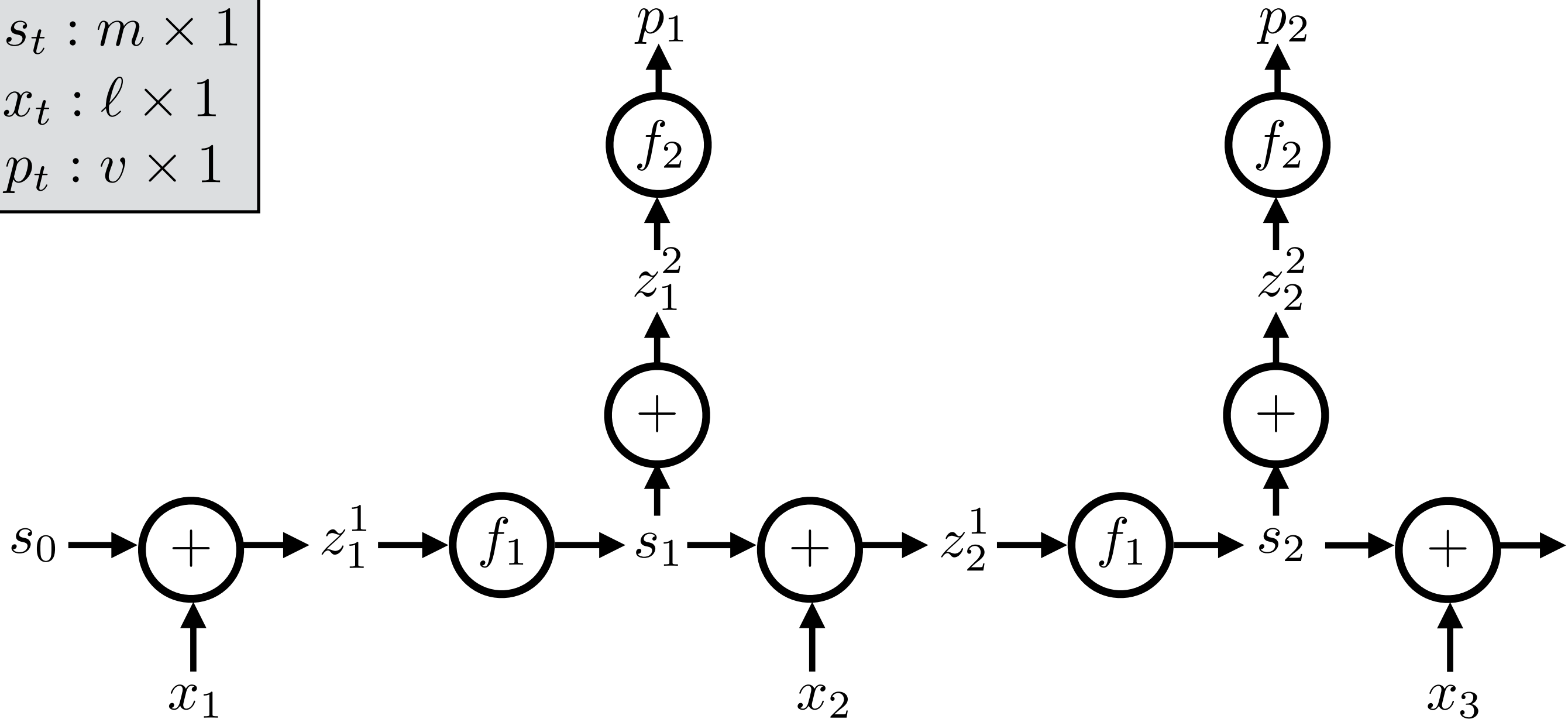
$$\begin{array}{l} s_t : m \times 1 \\ x_t : \ell \times 1 \\ p_t : v \times 1 \end{array}$$



- Example: Alphabet of ℓ chars; state is last c chars ($m = c\ell$)
- $$s_t = f_1 \left(\underbrace{W^{sx} x_t + W^{ss} s_{t-1} + W_0^{ss}}_{z_t^1} \right)$$
- $$p_t = f_2 \left(\underbrace{W^o s_t + W_0^o}_{z_t^2} \right)$$
- $$p^{(i)} = \text{RNN}(x^{(i)}; W, W_0)$$

Recurrent neural network

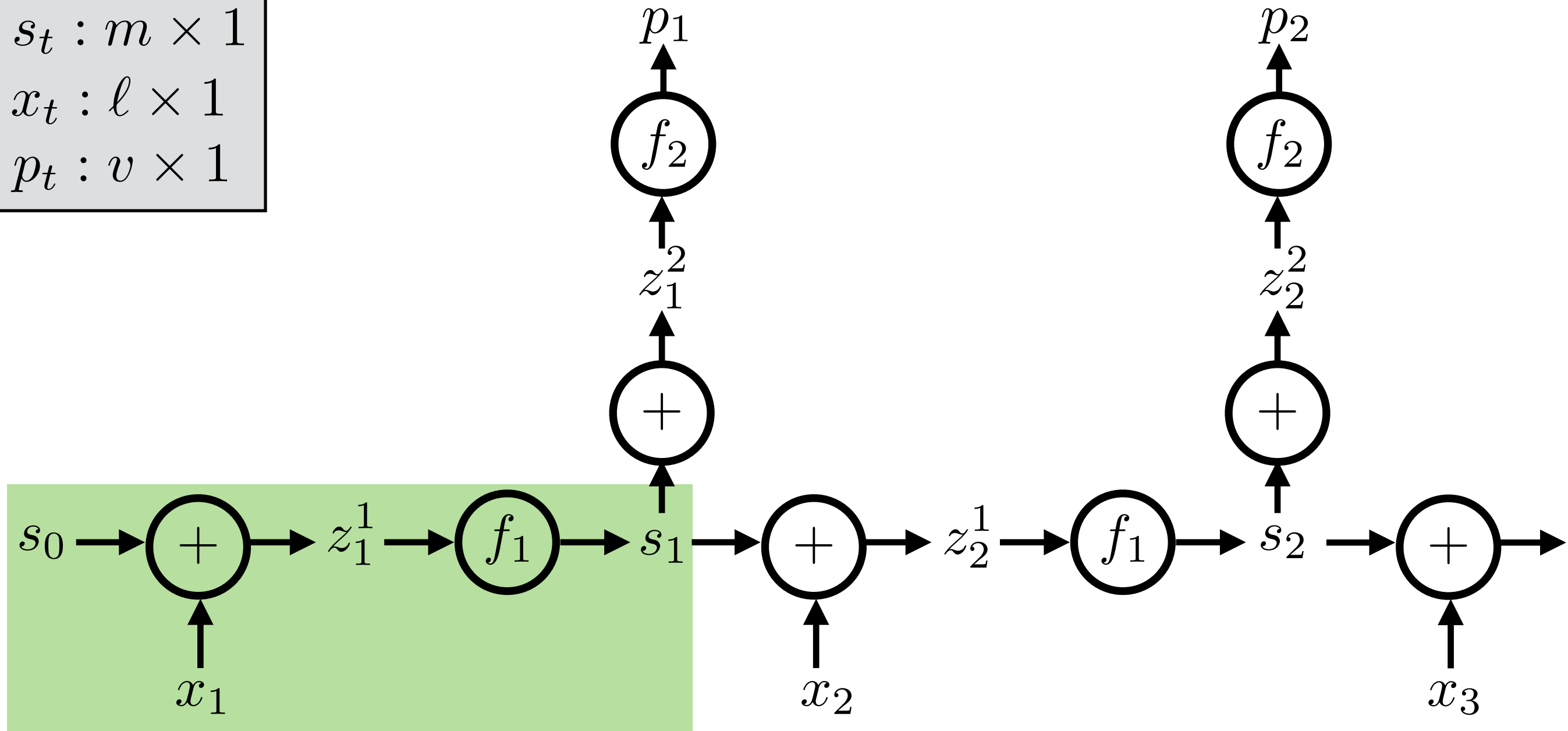
$$\begin{array}{l} s_t : m \times 1 \\ x_t : \ell \times 1 \\ p_t : v \times 1 \end{array}$$



- Example: Alphabet of ℓ chars; state is last c chars ($m = c\ell$)
- $$s_t = f_1 \left(\underbrace{W^{sx} x_t + W^{ss} s_{t-1} + W_0^{ss}}_{z_t^1} \right)$$
- $$p_t = f_2 \left(\underbrace{W^o s_t + W_0^o}_{z_t^2} \right)$$
- $$p^{(i)} = \text{RNN}(x^{(i)}; W, W_0)$$

Recurrent neural network

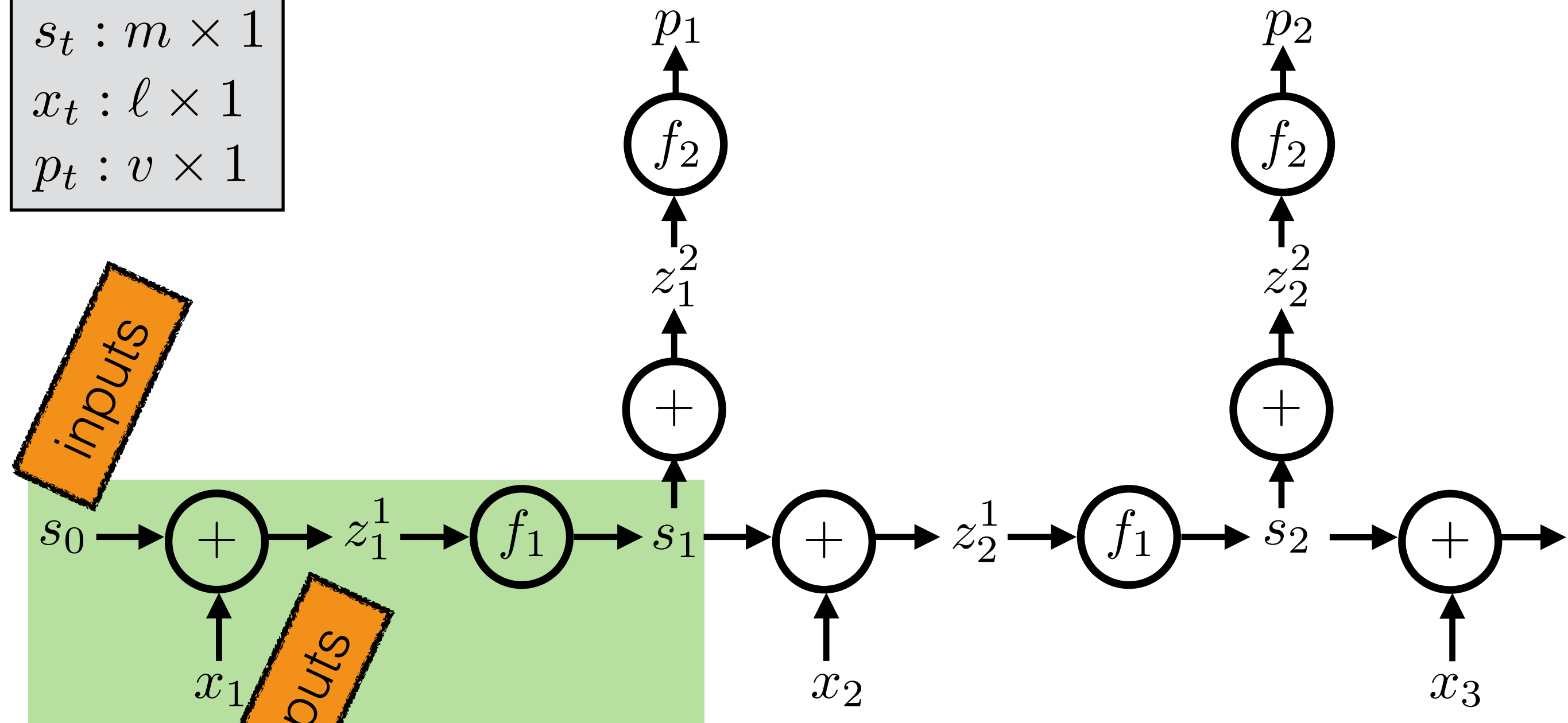
$$\begin{array}{l} s_t : m \times 1 \\ x_t : \ell \times 1 \\ p_t : v \times 1 \end{array}$$



- Example: Alphabet of ℓ chars; state is last c chars ($m = c\ell$)
- $$s_t = f_1 \left(\underbrace{W^{sx} x_t + W^{ss} s_{t-1} + W_0^{ss}}_{z_t^1} \right)$$
- $$p_t = f_2 \left(\underbrace{W^o s_t + W_0^o}_{z_t^2} \right)$$
- $$p^{(i)} = \text{RNN}(x^{(i)}; W, W_0)$$

Recurrent neural network

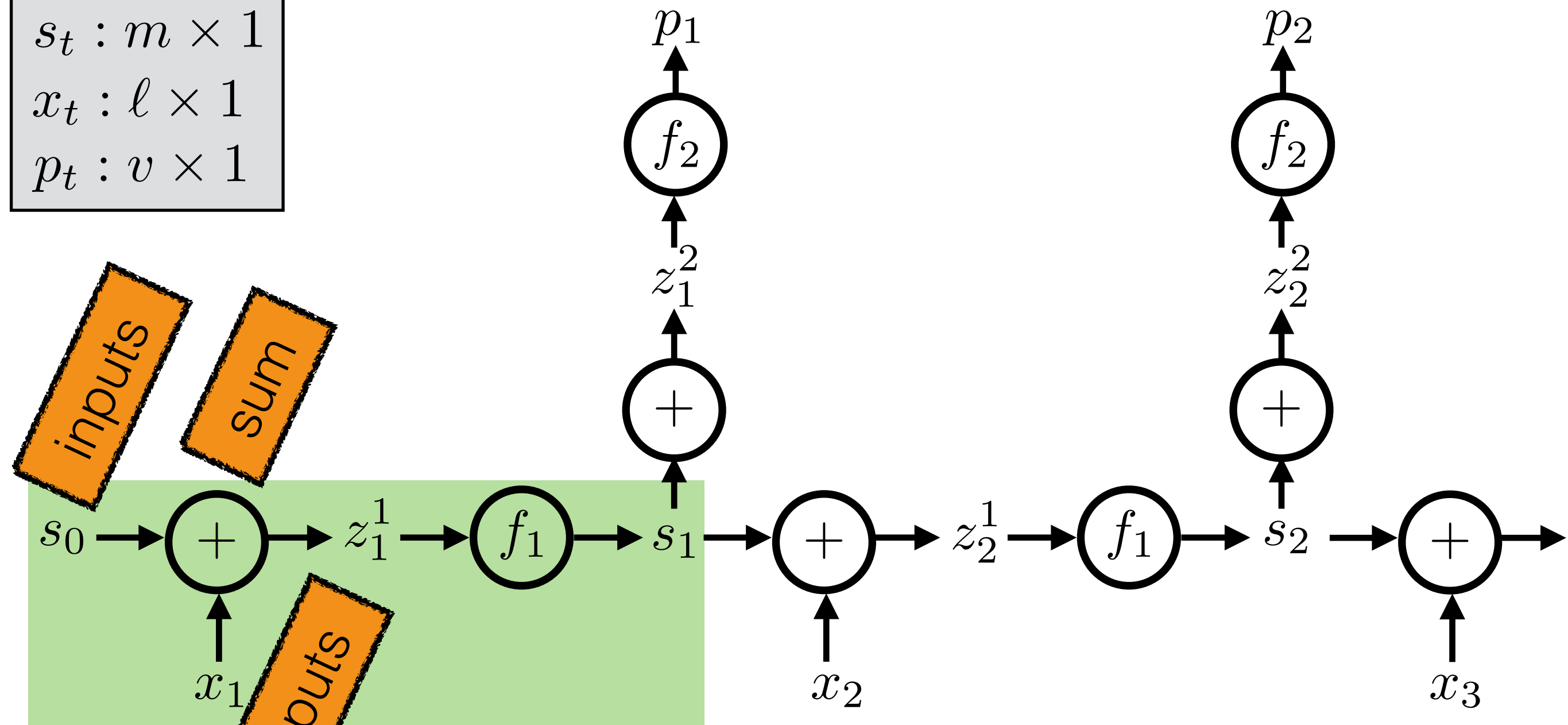
$$\begin{aligned} s_t &: m \times 1 \\ x_t &: \ell \times 1 \\ p_t &: v \times 1 \end{aligned}$$



- Example: Alphabet of ℓ chars; state is last c chars ($m = c\ell$)
- $$s_t = f_1 \left(\underbrace{W^{sx} x_t + W^{ss} s_{t-1} + W_0^{ss}}_{z_t^1} \right)$$
- $$p_t = f_2 \left(\underbrace{W^o s_t + W_0^o}_{z_t^2} \right)$$
- $$p^{(i)} = \text{RNN}(x^{(i)}; W, W_0)$$

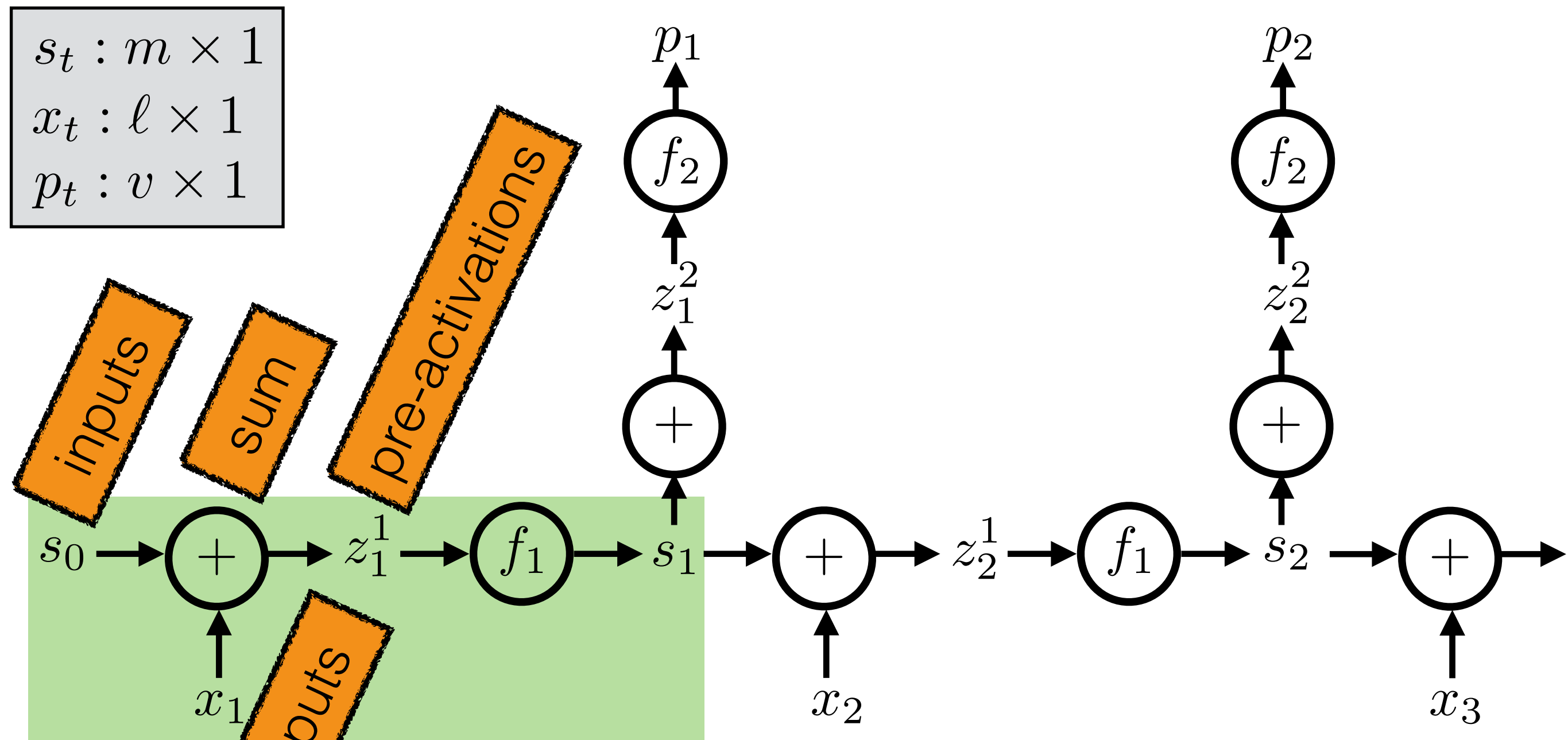
Recurrent neural network

$$\begin{aligned} s_t &: m \times 1 \\ x_t &: \ell \times 1 \\ p_t &: v \times 1 \end{aligned}$$



- Example: Alphabet of ℓ chars; state is last c chars ($m = c\ell$)
- $$s_t = f_1 \left(\underbrace{W^{sx} x_t + W^{ss} s_{t-1} + W_0^{ss}}_{z_t^1} \right)$$
- $$p_t = f_2 \left(\underbrace{W^o s_t + W_0^o}_{z_t^2} \right)$$
- $$p^{(i)} = \text{RNN}(x^{(i)}; W, W_0)$$

Recurrent neural network



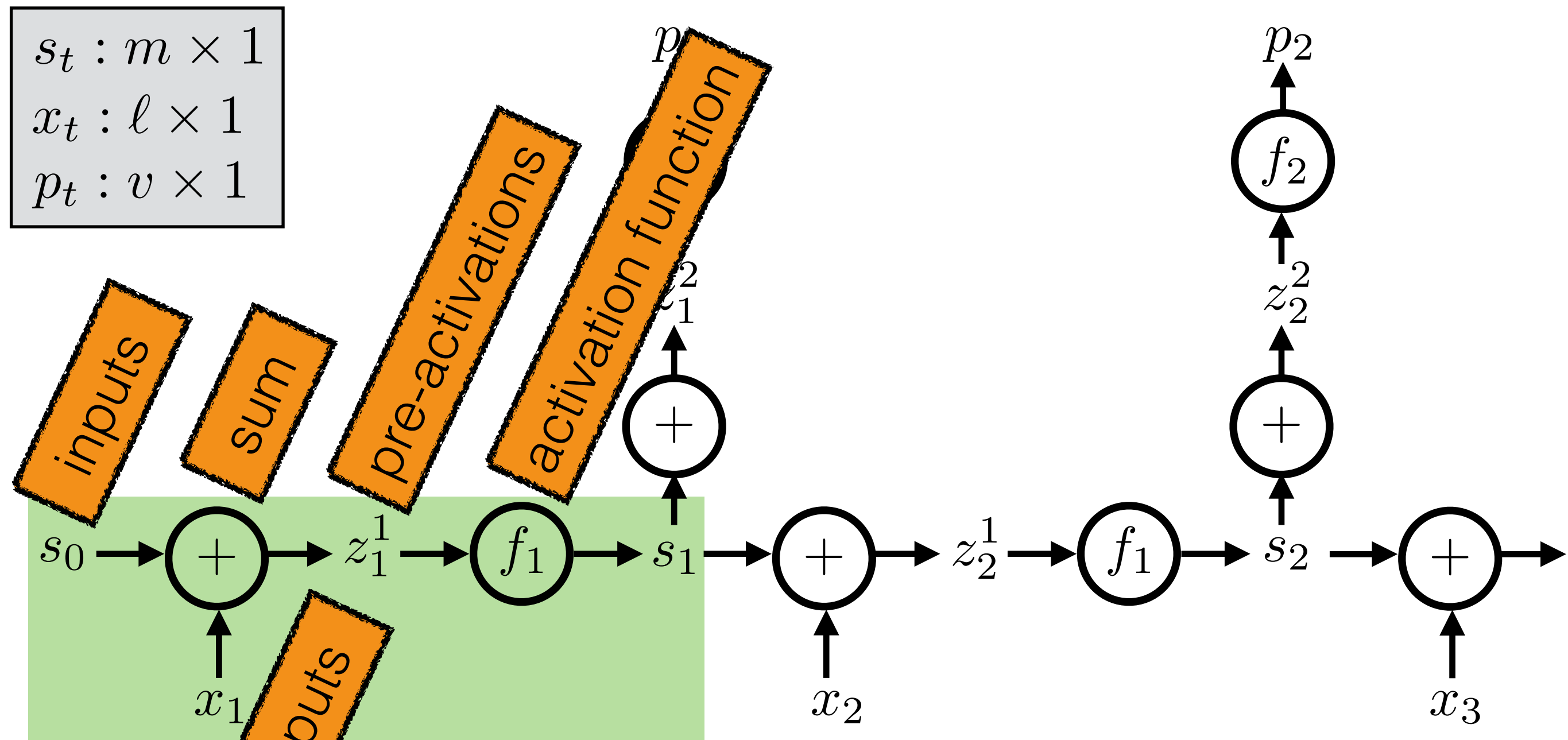
- Example: Alphabet of ℓ chars; state is last c chars ($m = c\ell$)

$$s_t = f_1 \left(\underbrace{W^{sx} x_t + W^{ss} s_{t-1} + W_0^{ss}}_{z_t^1} \right)$$

$$p_t = f_2 \left(\underbrace{W^o s_t + W_0^o}_{z_t^2} \right)$$

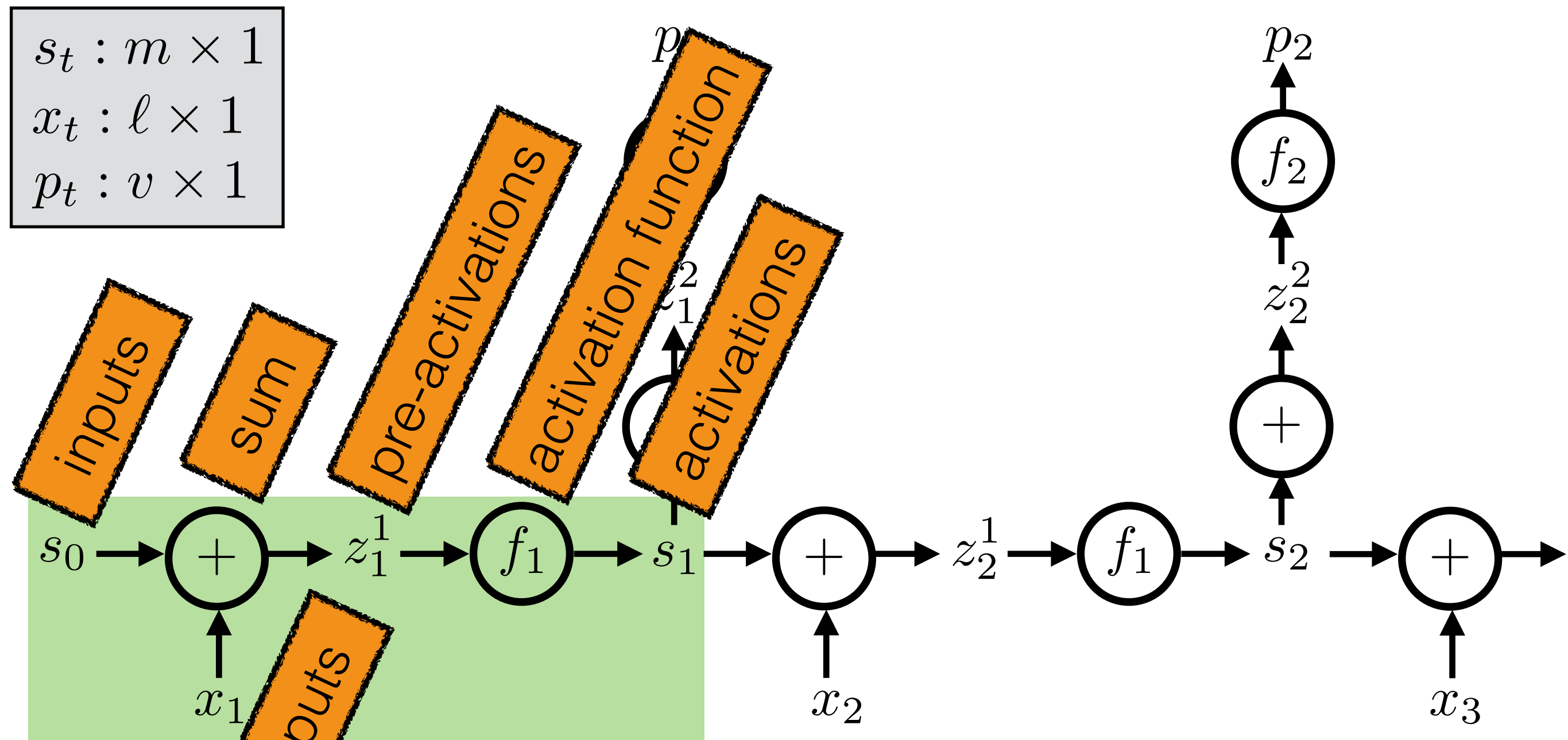
$$p^{(i)} = \text{RNN}(x^{(i)}; W, W_0)$$

Recurrent neural network



- Example: Alphabet of ℓ chars; state is last c chars ($m = c\ell$)
- $$s_t = f_1 \left(\underbrace{W^{sx} x_t + W^{ss} s_{t-1} + W_0^{ss}}_{z_t^1} \right)$$
- $$p_t = f_2 \left(\underbrace{W^o s_t + W_0^o}_{z_t^2} \right)$$
- $$p^{(i)} = \text{RNN}(x^{(i)}; W, W_0)$$

Recurrent neural network



- Example: Alphabet of ℓ chars; state is last c chars ($m = c\ell$)

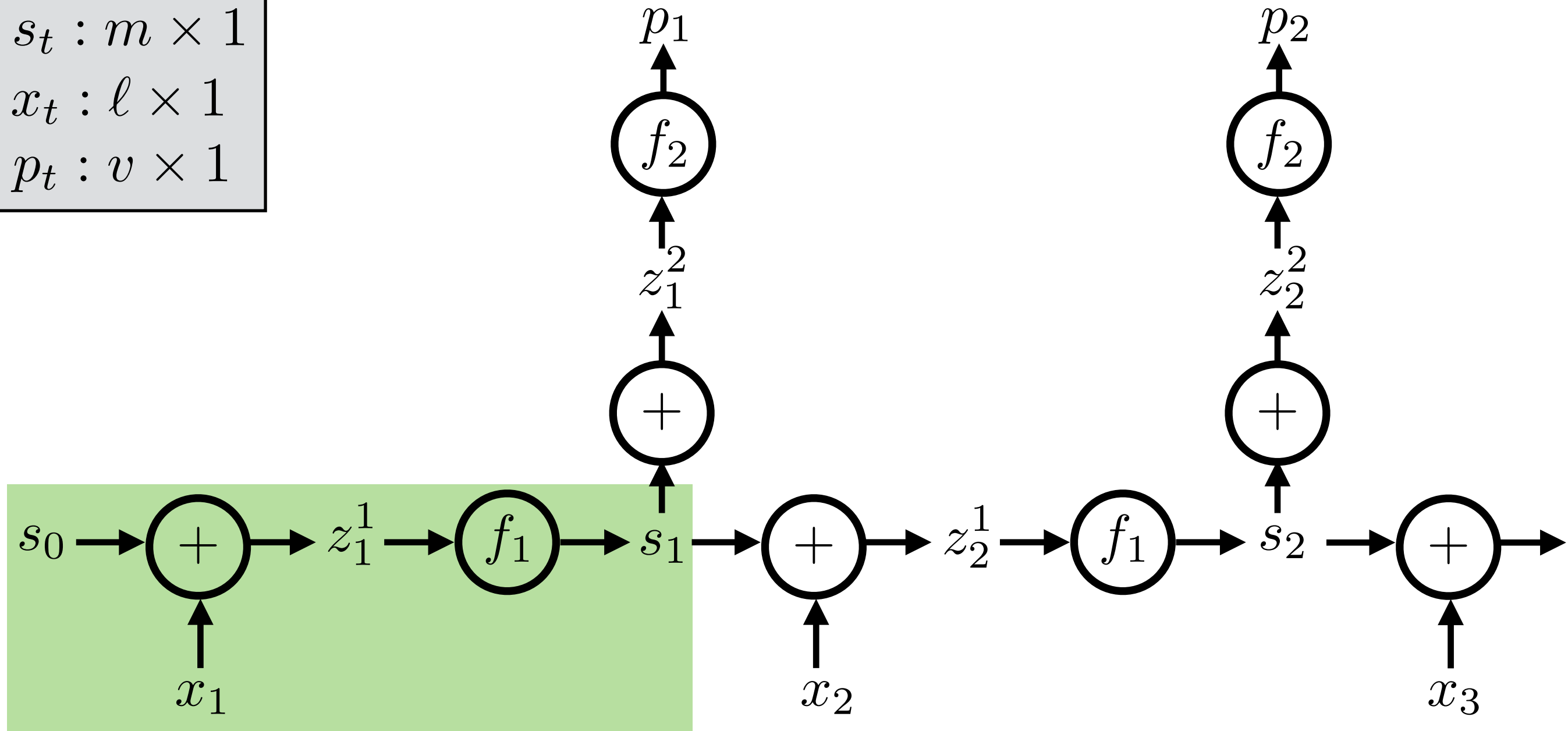
$$s_t = f_1 \left(\underbrace{W^{sx} x_t + W^{ss} s_{t-1} + W_0^{ss}}_{z_t^1} \right)$$

$$p_t = f_2 \left(\underbrace{W^o s_t + W_0^o}_{z_t^2} \right)$$

$$p^{(i)} = \text{RNN}(x^{(i)}; W, W_0)$$

Recurrent neural network

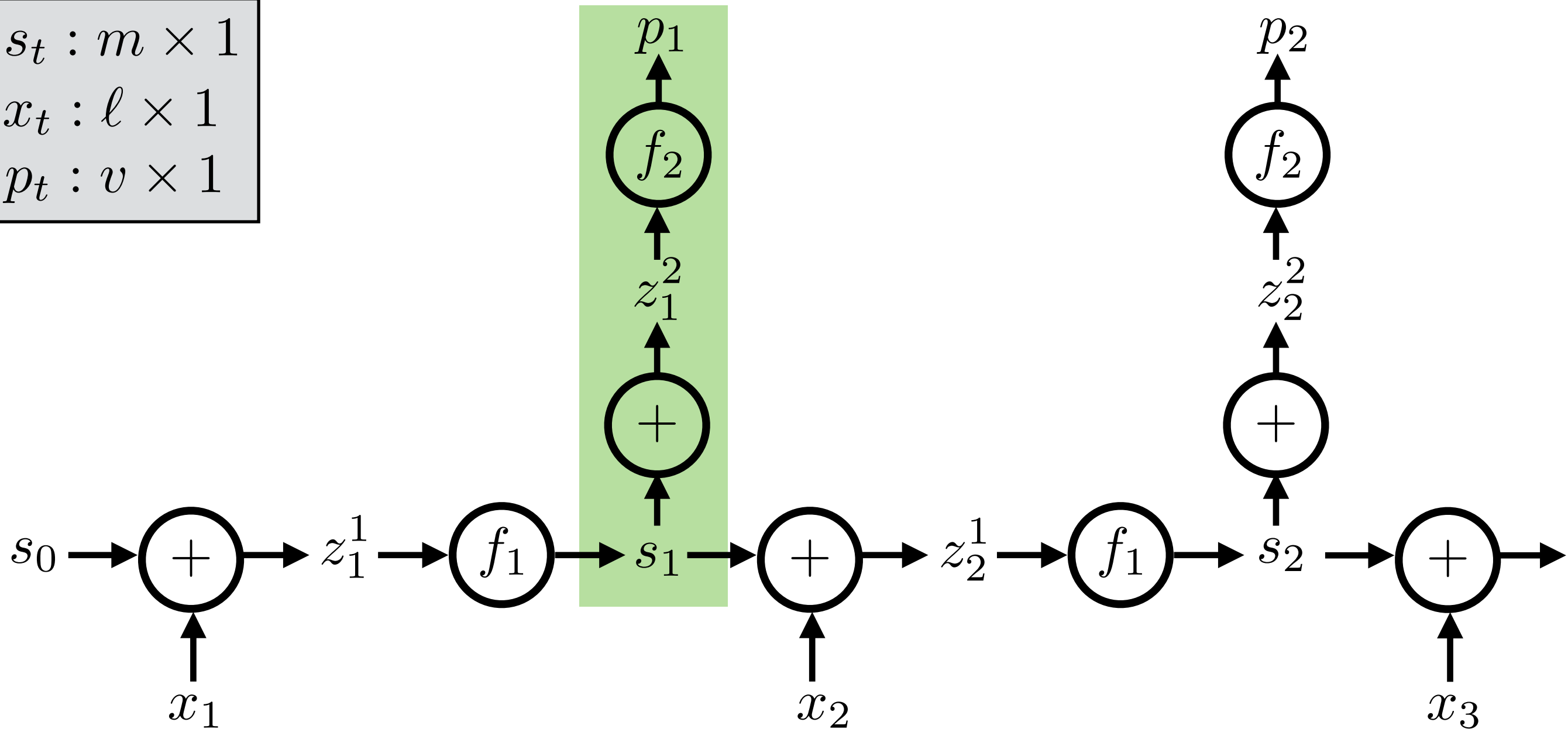
$$\begin{array}{l} s_t : m \times 1 \\ x_t : \ell \times 1 \\ p_t : v \times 1 \end{array}$$



- Example: Alphabet of ℓ chars; state is last c chars ($m = c\ell$)
- $$s_t = f_1 \left(\underbrace{W^{sx} x_t + W^{ss} s_{t-1} + W_0^{ss}}_{z_t^1} \right)$$
- $$p_t = f_2 \left(\underbrace{W^o s_t + W_0^o}_{z_t^2} \right)$$
- $$p^{(i)} = \text{RNN}(x^{(i)}; W, W_0)$$

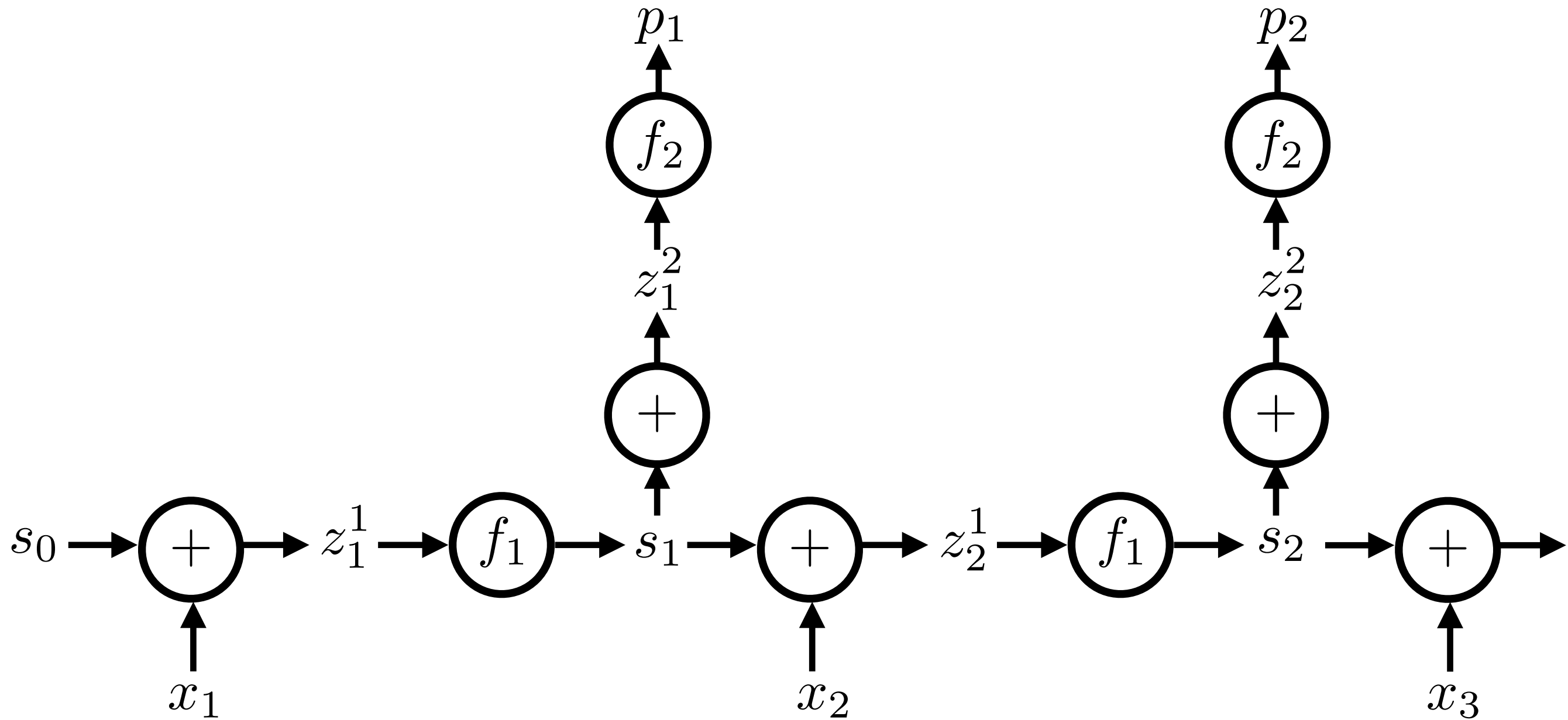
Recurrent neural network

$$\begin{array}{l} s_t : m \times 1 \\ x_t : \ell \times 1 \\ p_t : v \times 1 \end{array}$$

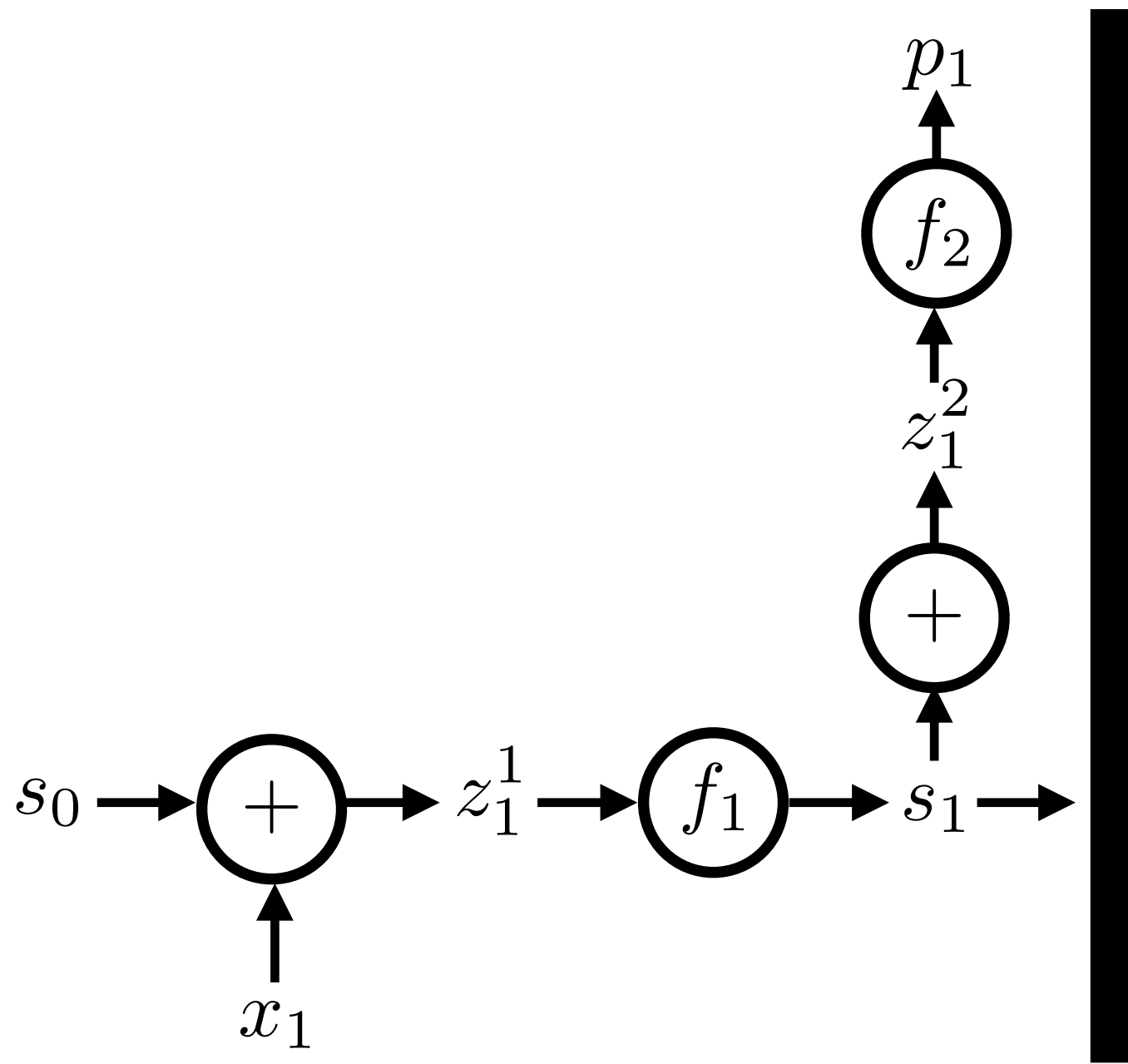


- Example: Alphabet of ℓ chars; state is last c chars ($m = c\ell$)
- $$s_t = f_1 \left(\underbrace{W^{sx} x_t + W^{ss} s_{t-1} + W_0^{ss}}_{z_t^1} \right)$$
- $$p_t = f_2 \left(\underbrace{W^o s_t + W_0^o}_{z_t^2} \right)$$
- $$p^{(i)} = \text{RNN}(x^{(i)}; W, W_0)$$

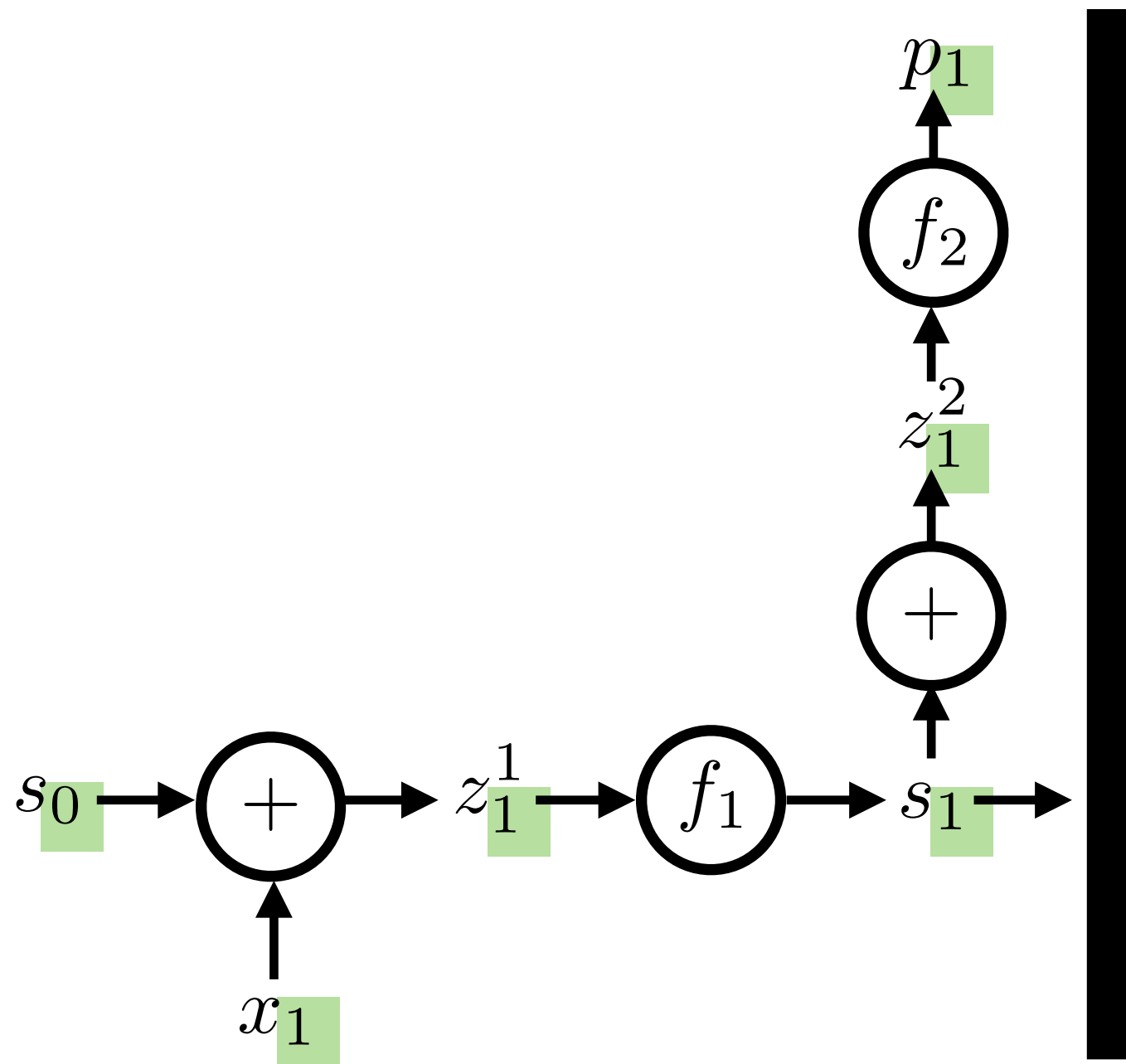
Recurrent neural network



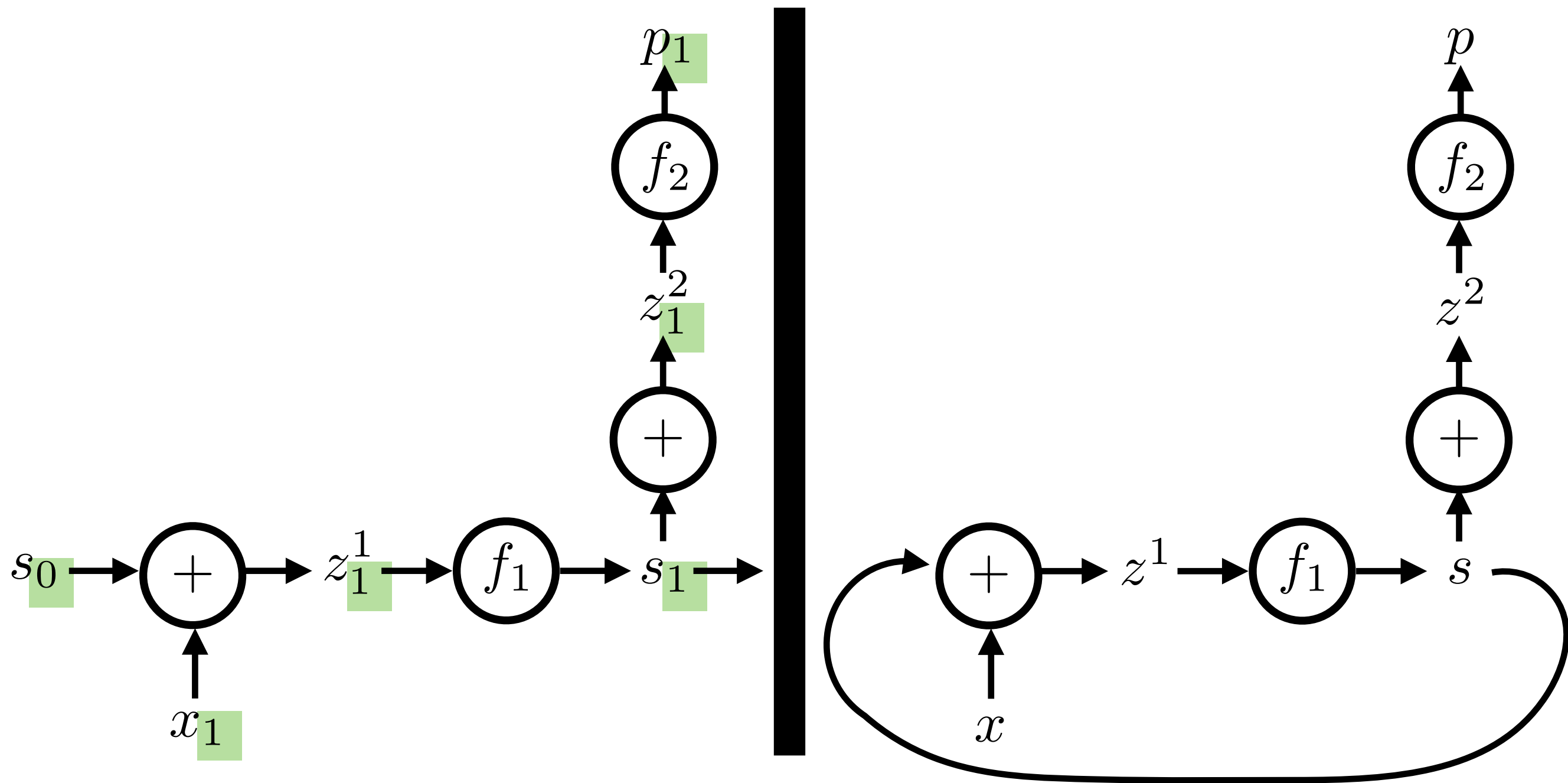
Recurrent neural network



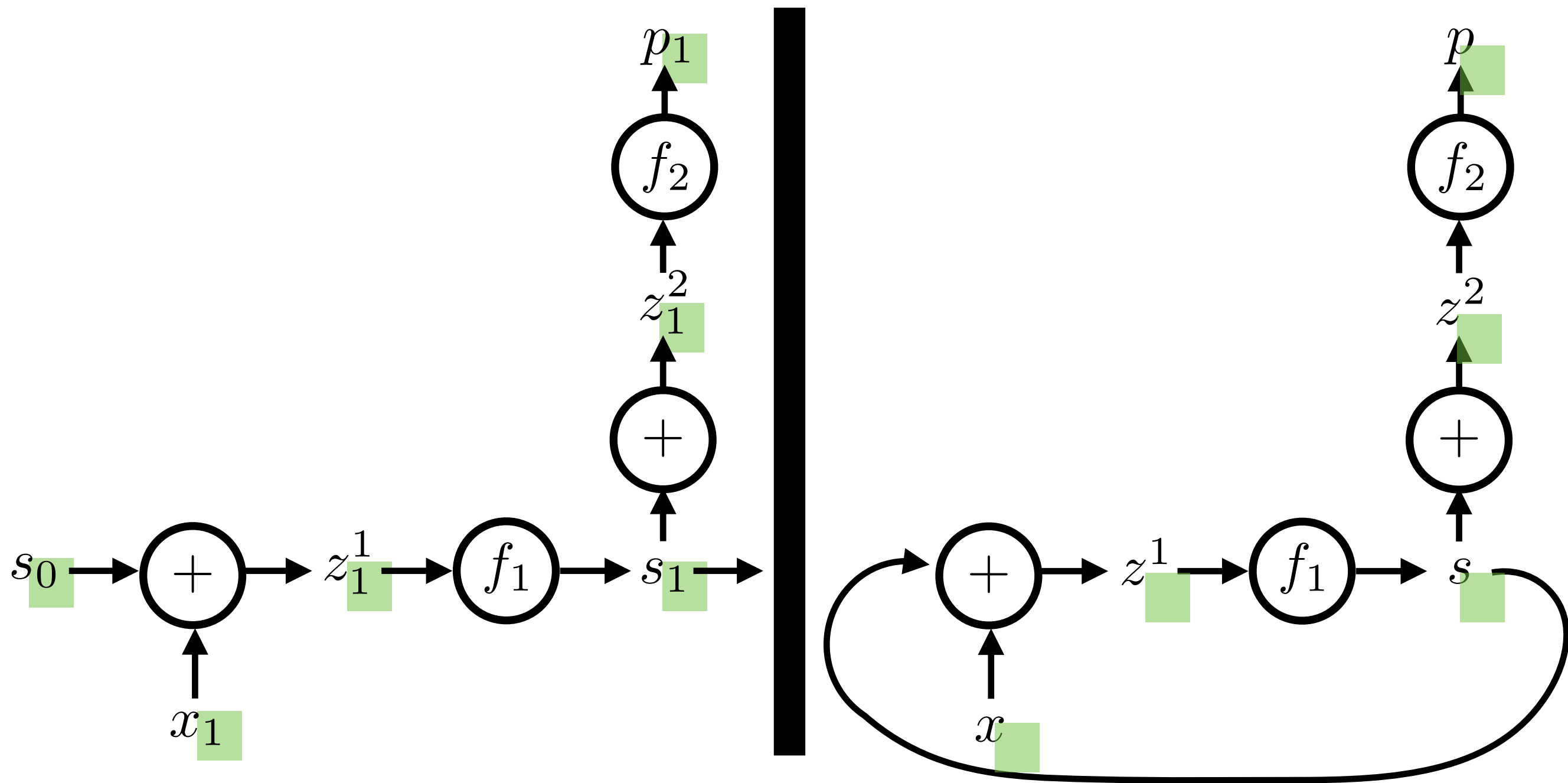
Recurrent neural network



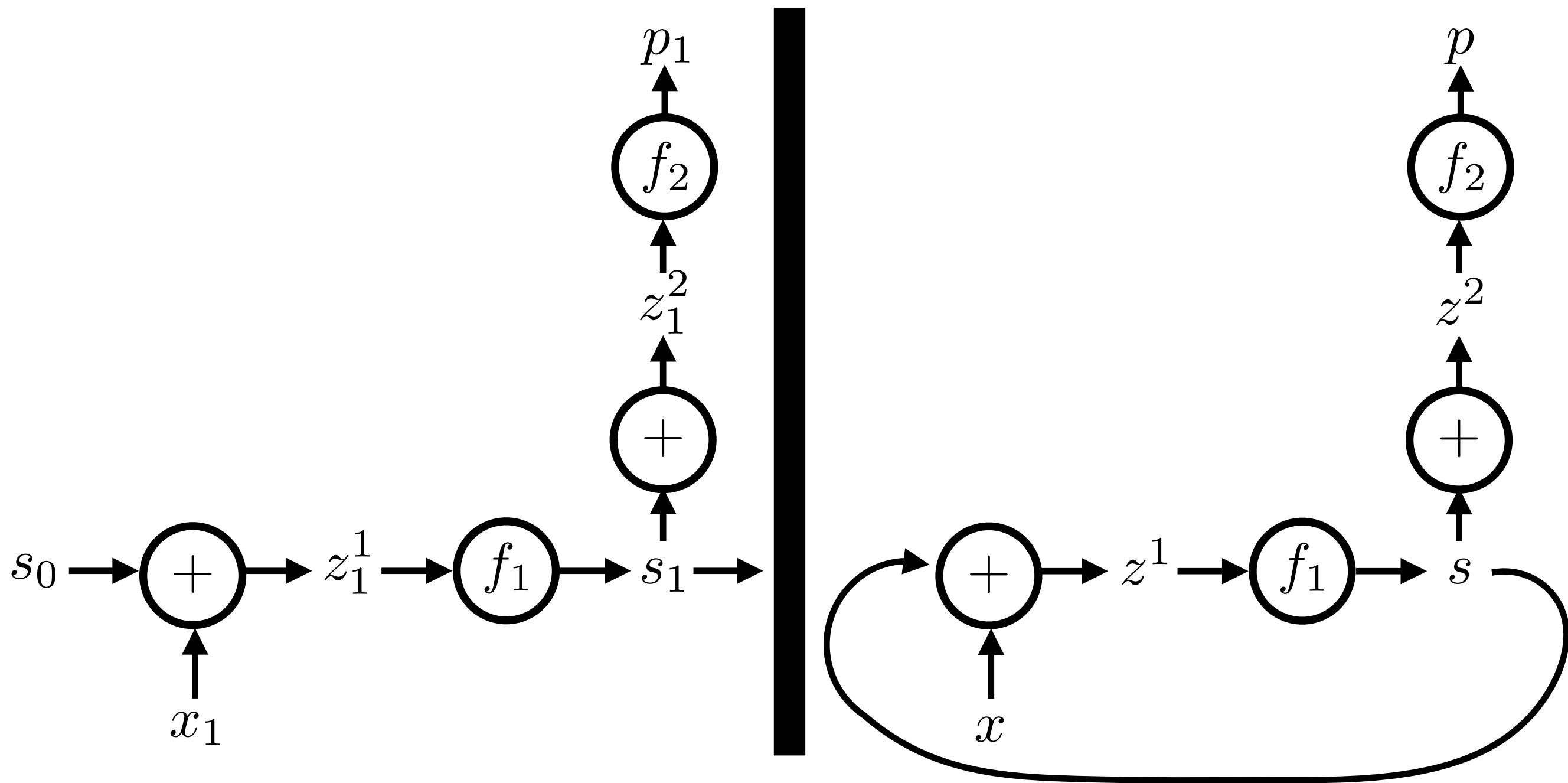
Recurrent neural network



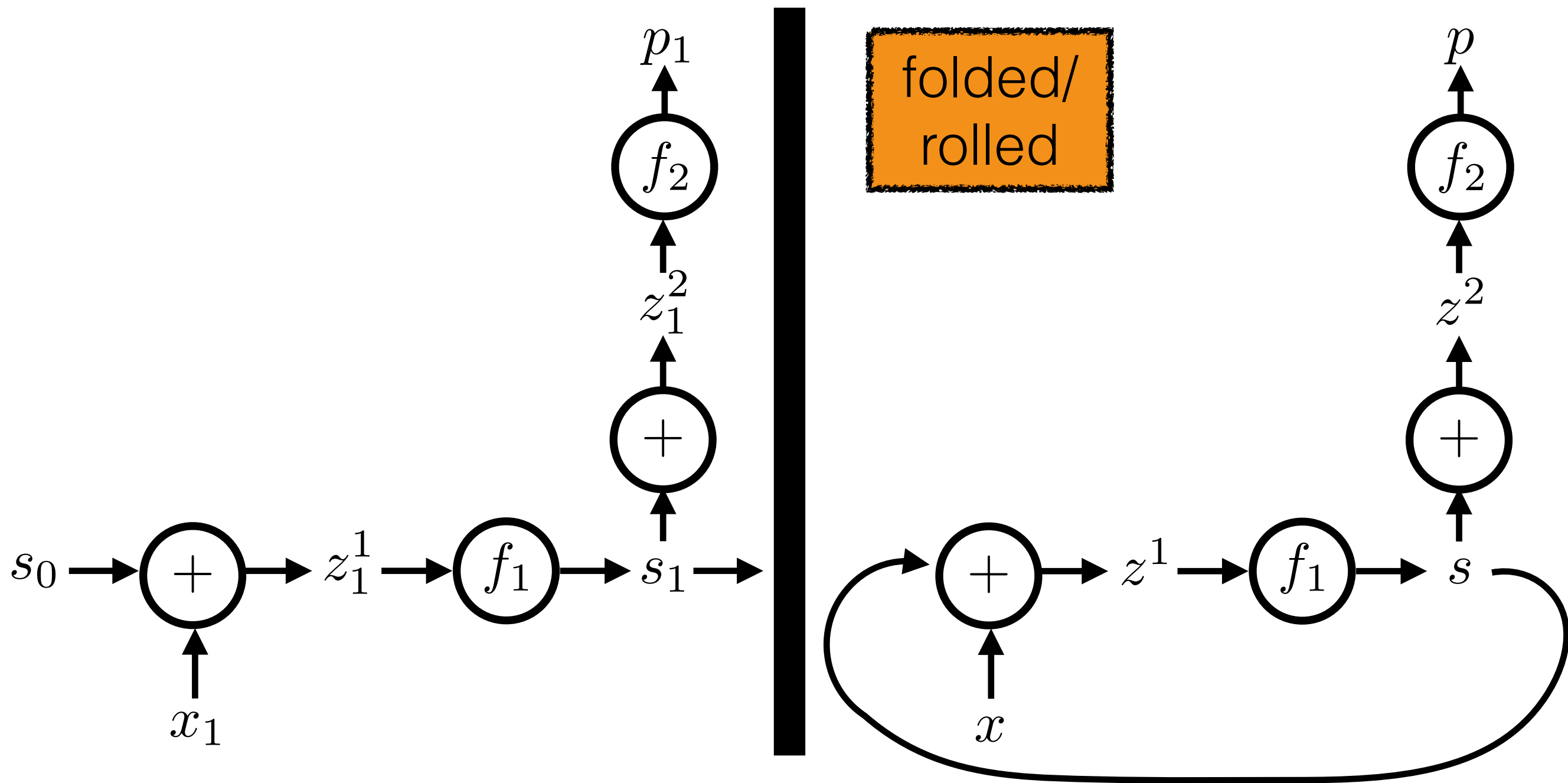
Recurrent neural network



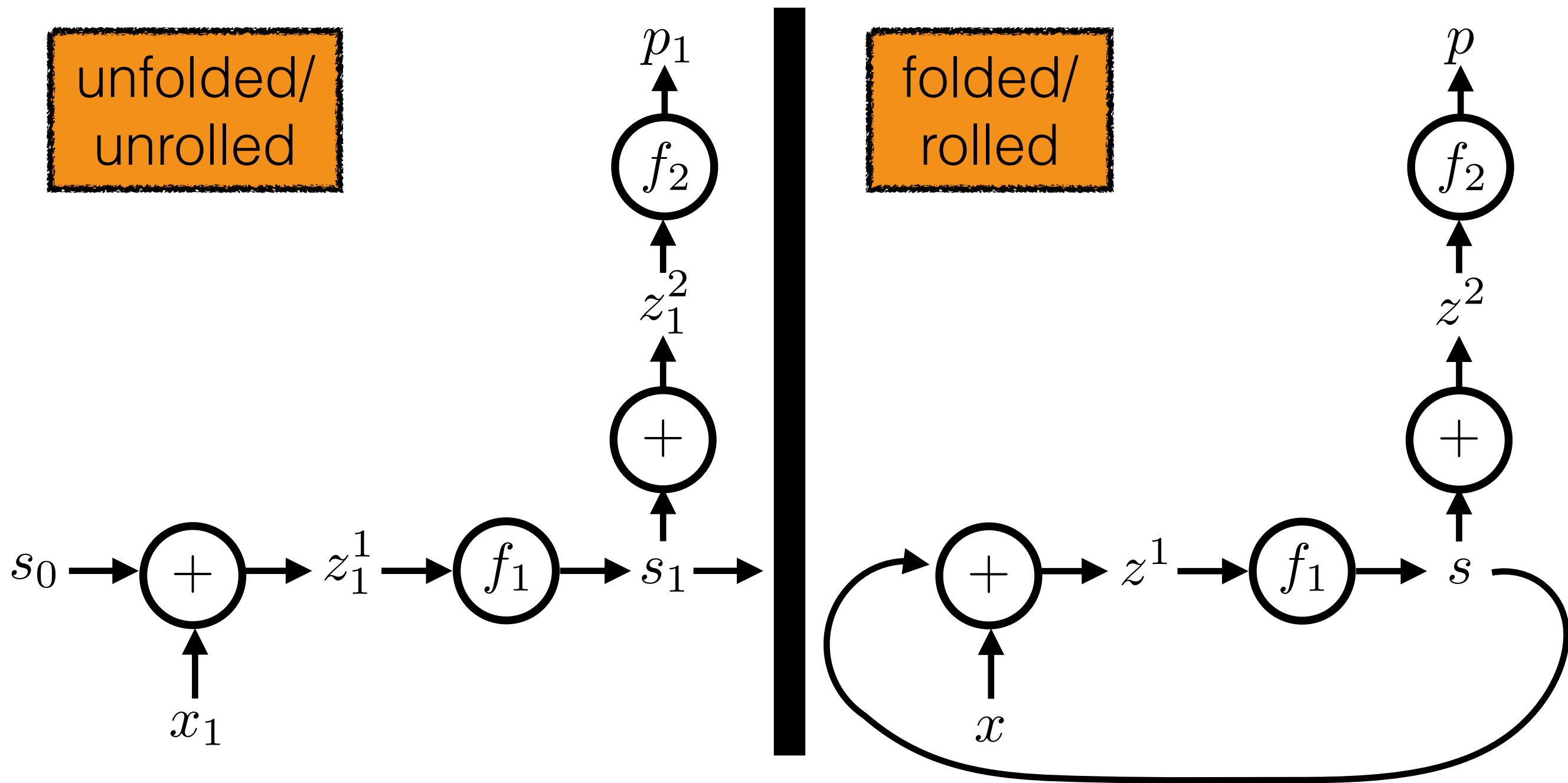
Recurrent neural network



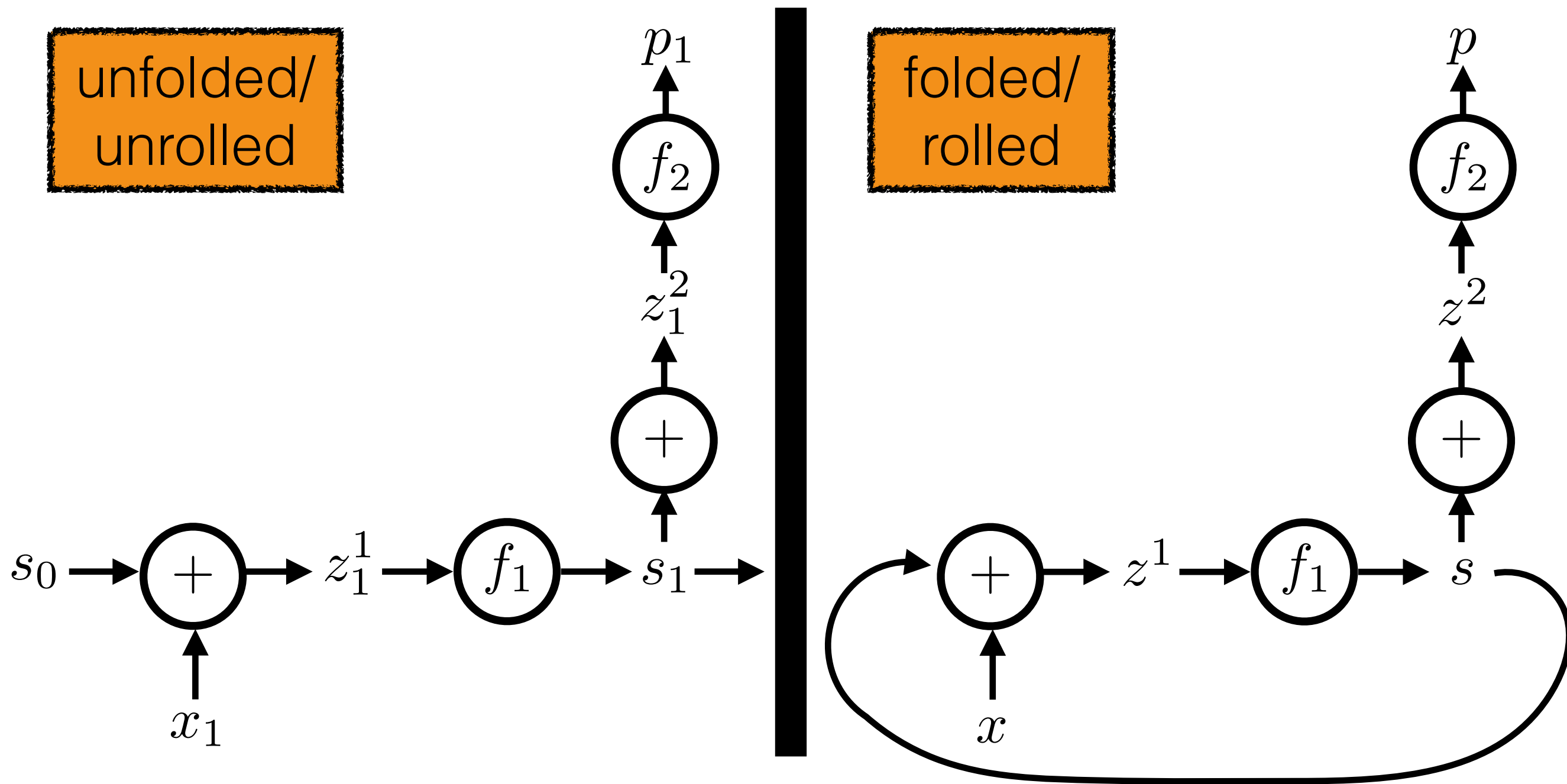
Recurrent neural network



Recurrent neural network

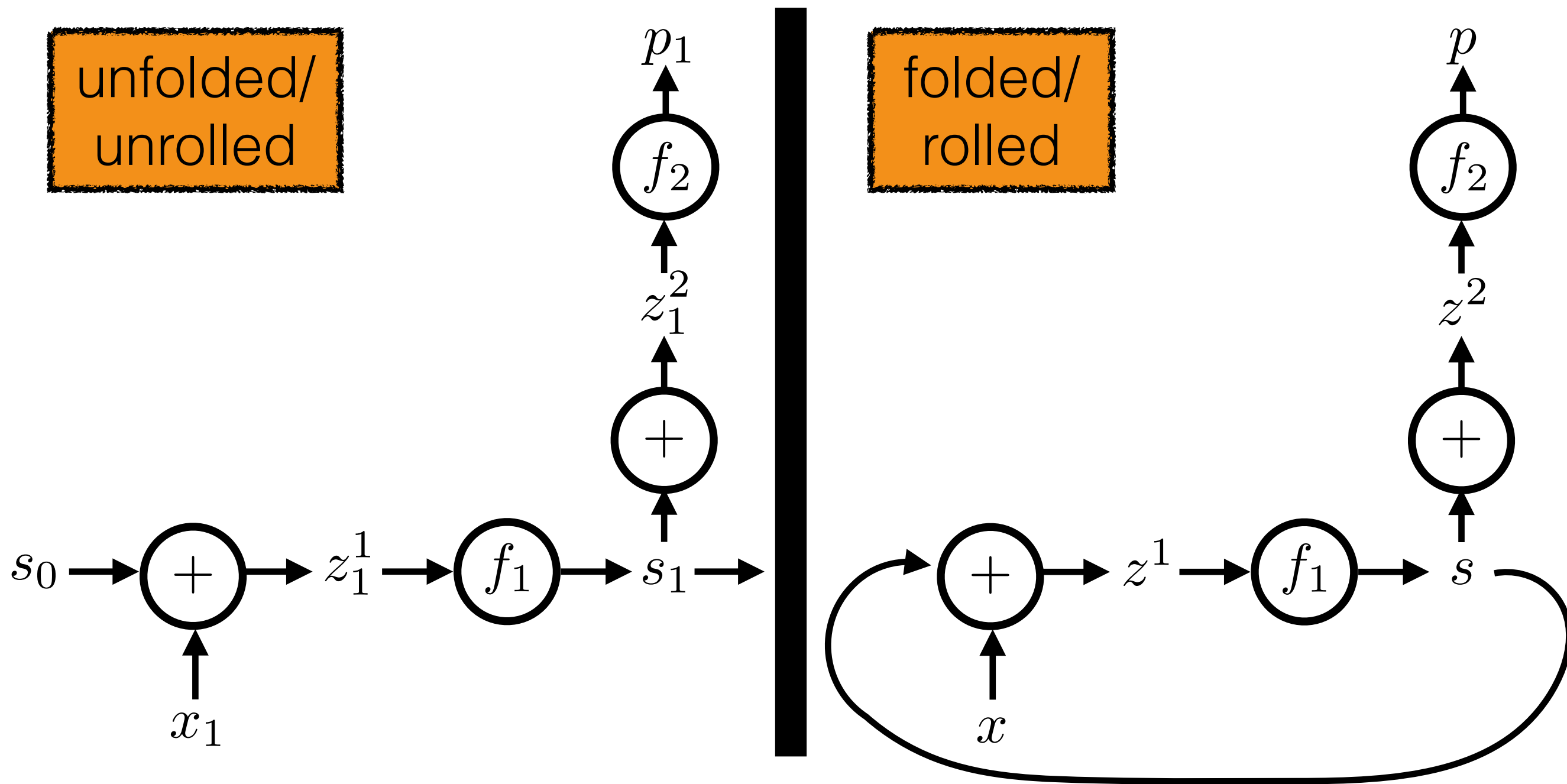


Recurrent neural network



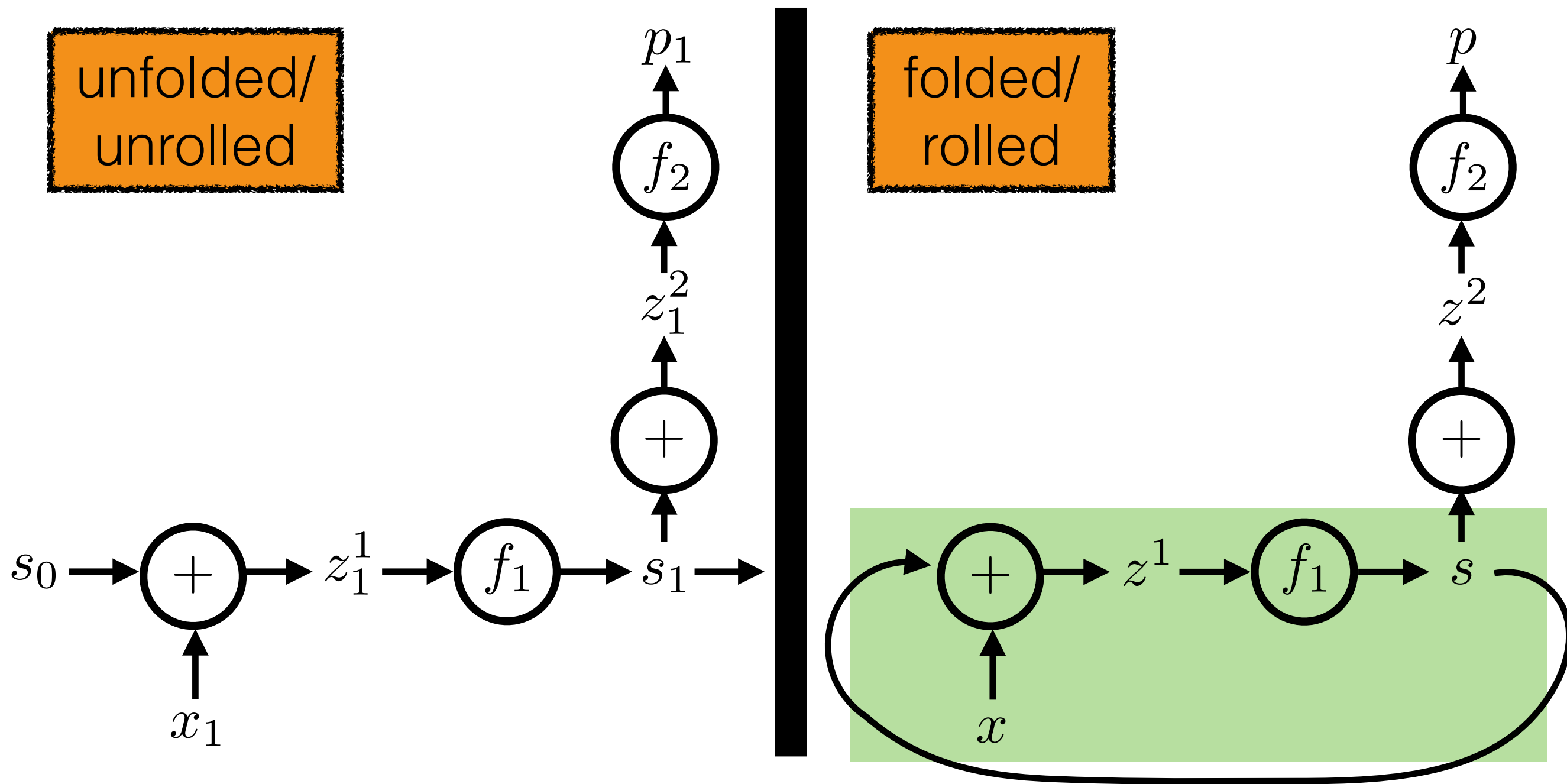
- Compare to:

Recurrent neural network



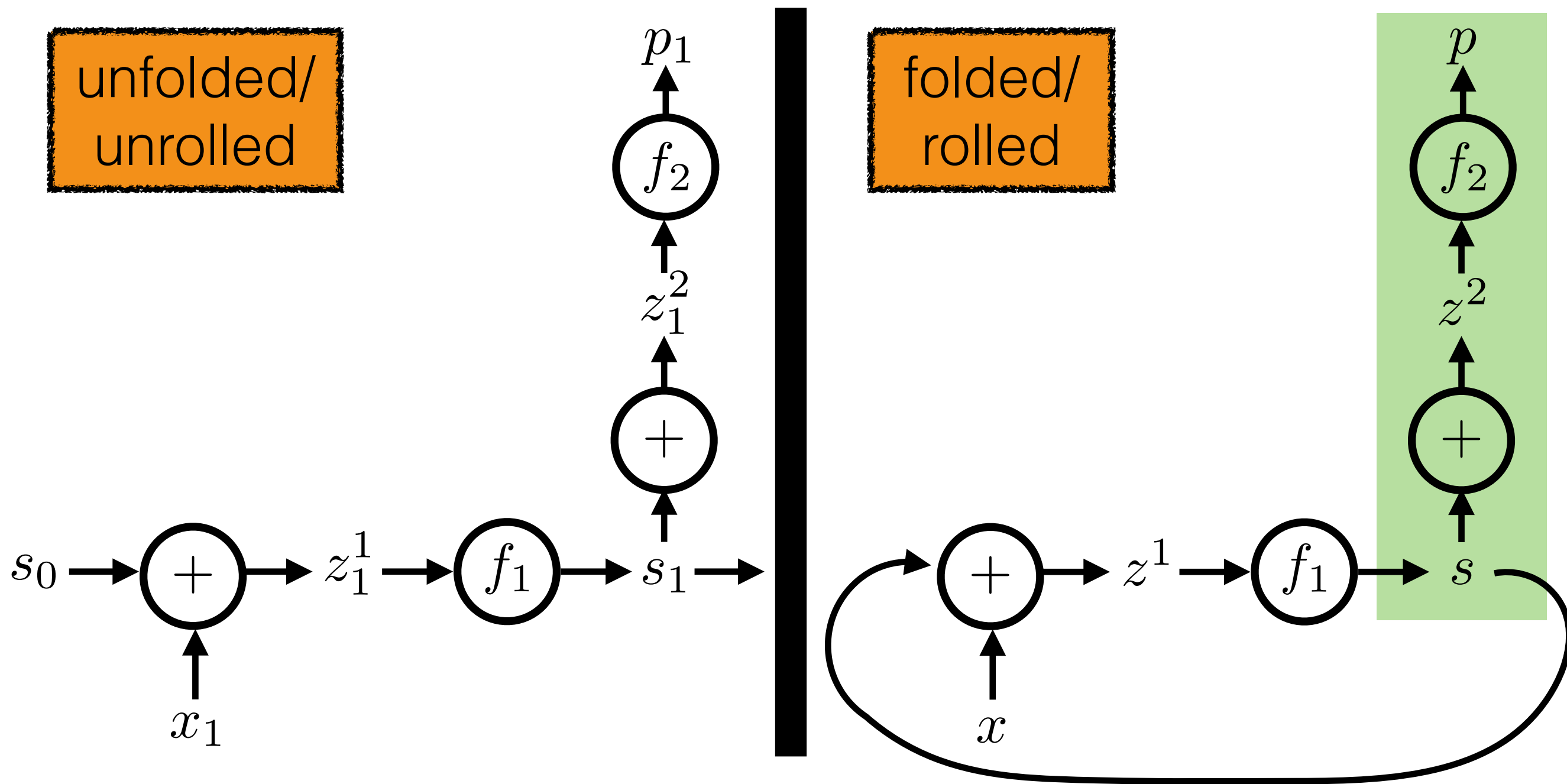
- Compare to:
 - Feedforward neural networks

Recurrent neural network



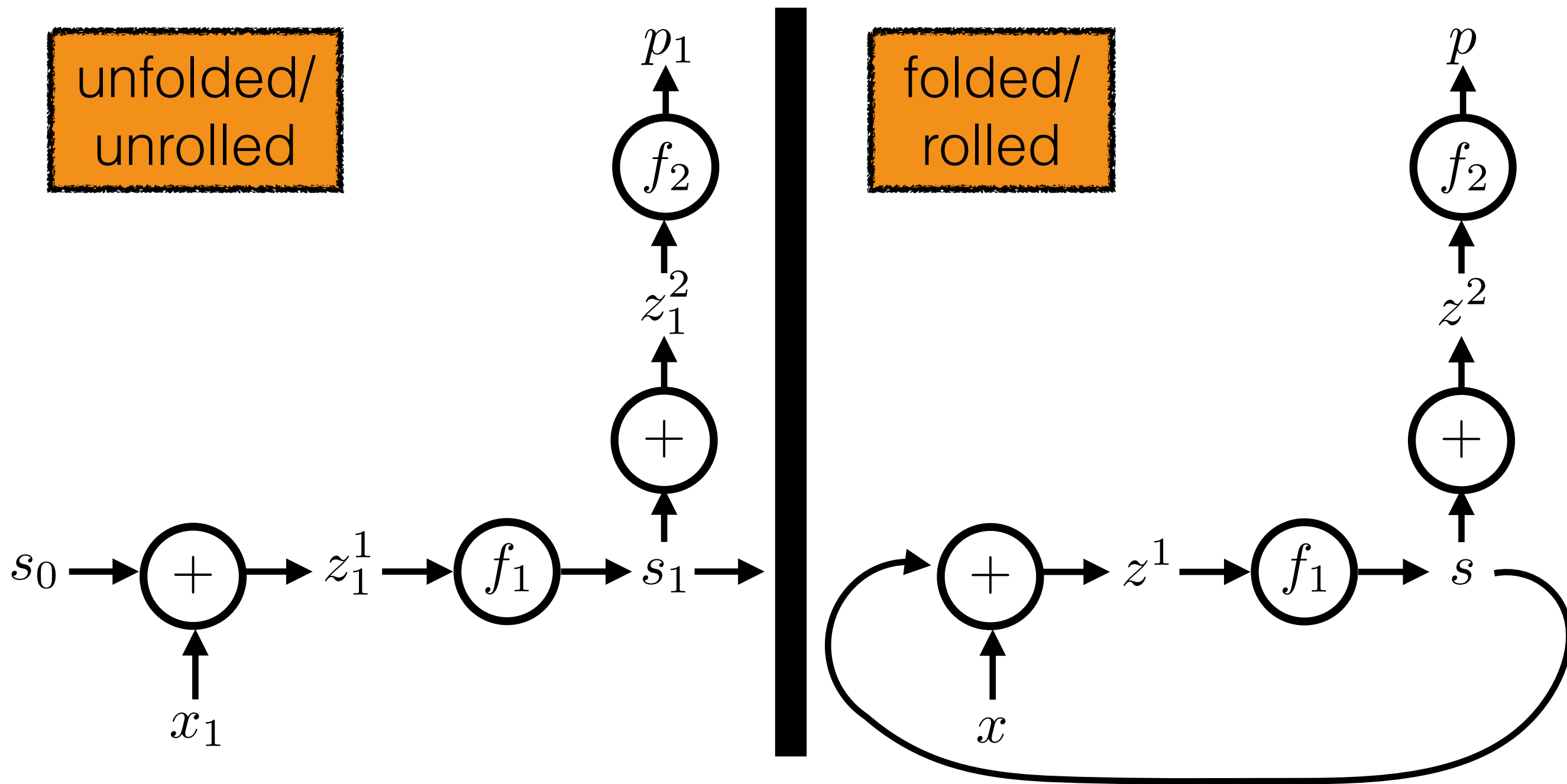
- Compare to:
 - Feedforward neural networks

Recurrent neural network



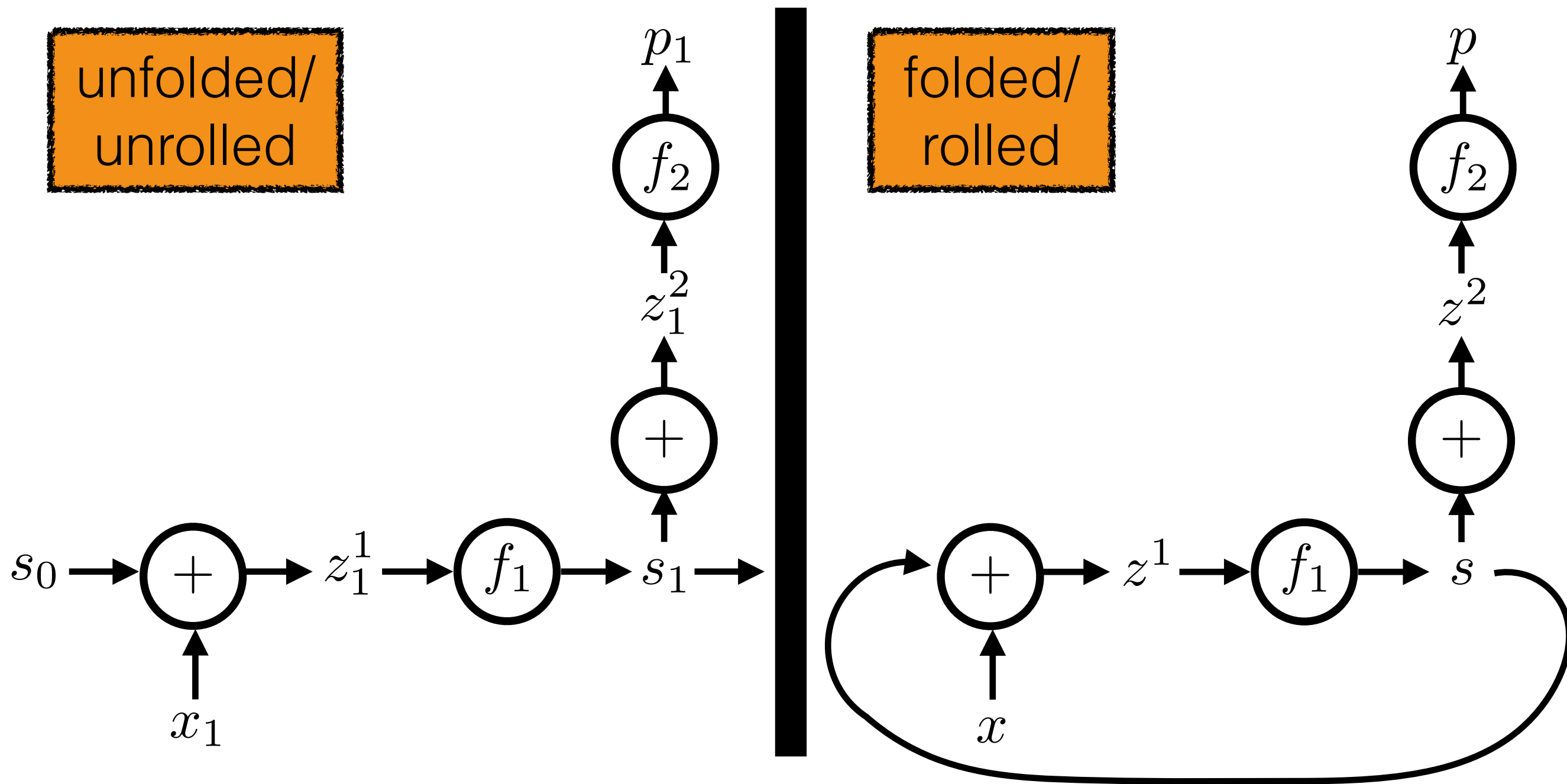
- Compare to:
 - Feedforward neural networks

Recurrent neural network

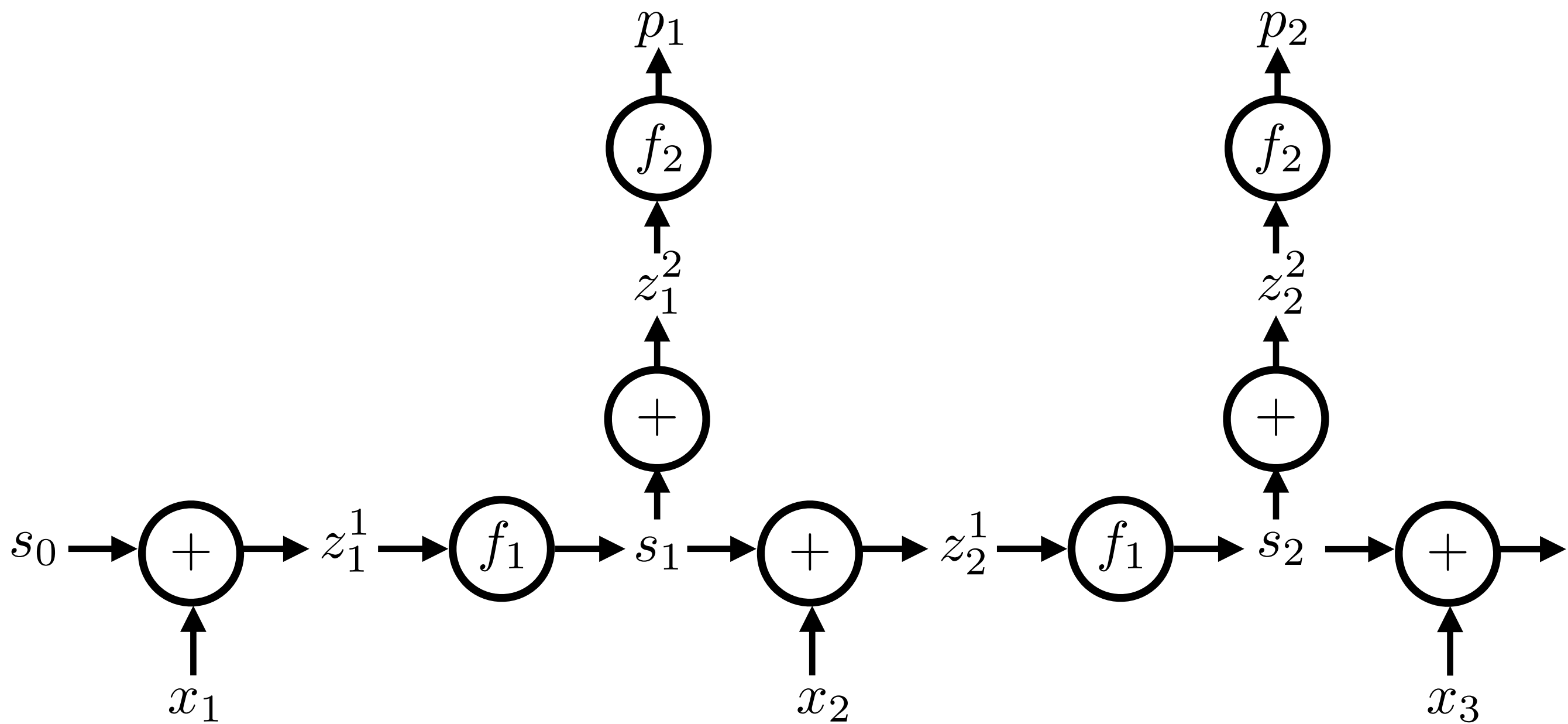


- Compare to:
 - Feedforward neural networks
 - Convolutional neural networks

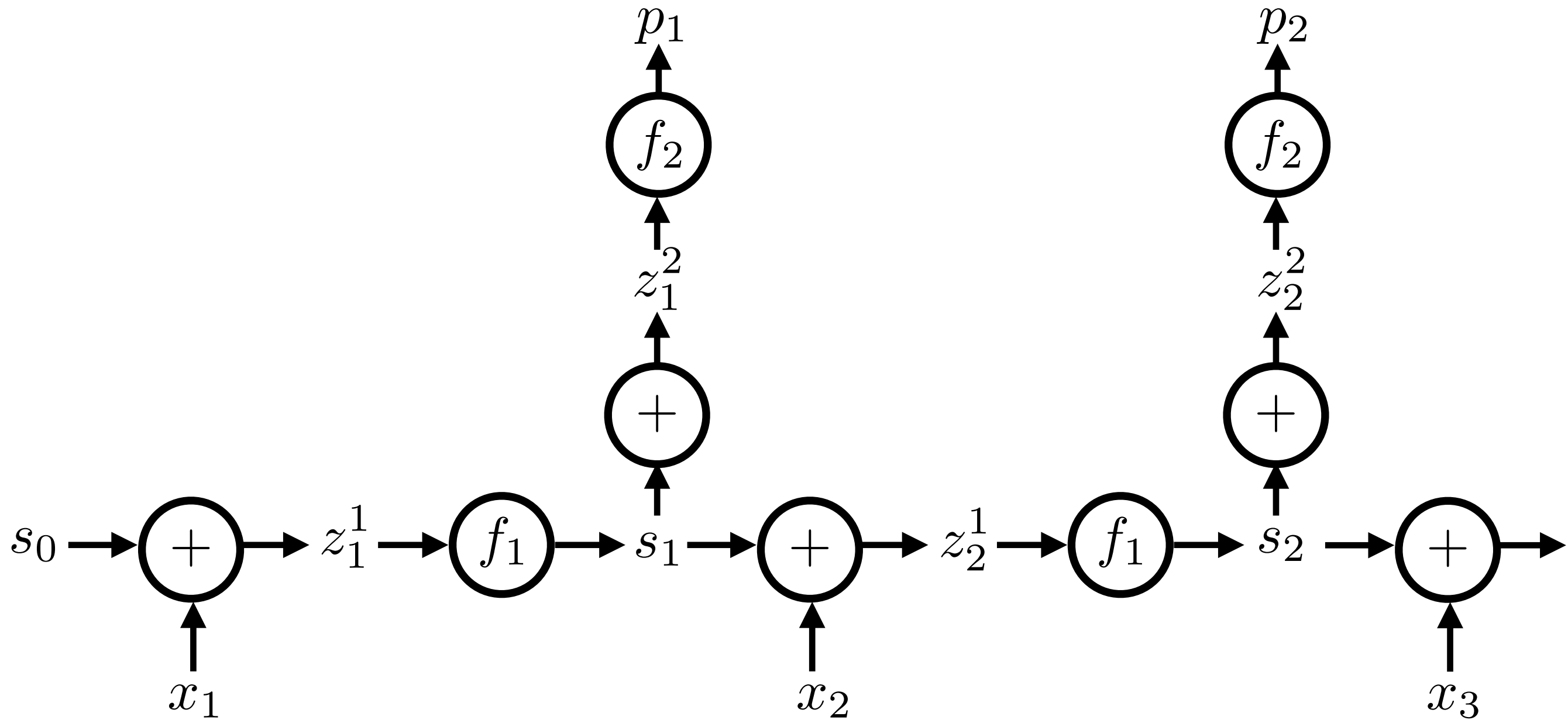
Recurrent neural network



- Compare to:
 - Feedforward neural networks
 - Convolutional neural networks
 - Reinforcement learning

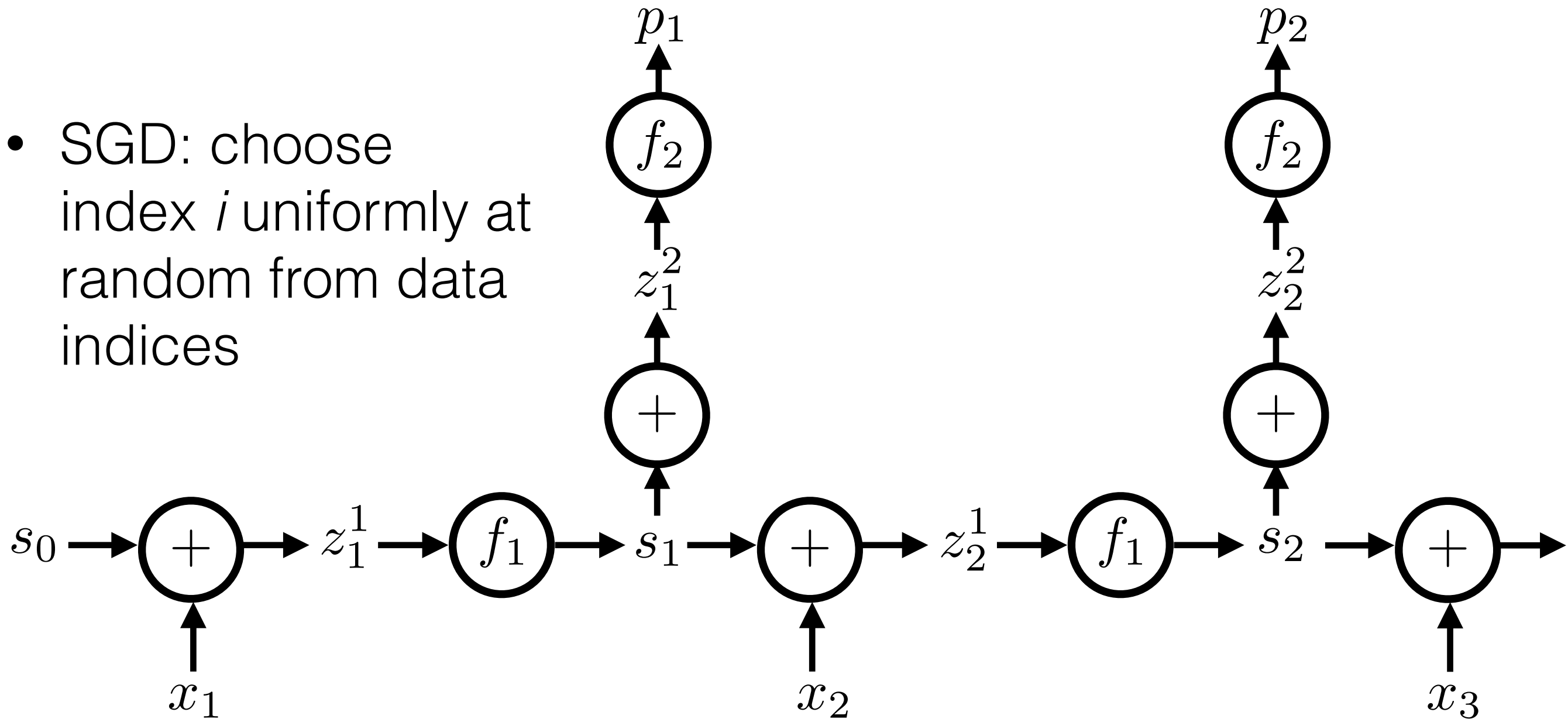


RNNs: a taste of backpropagation



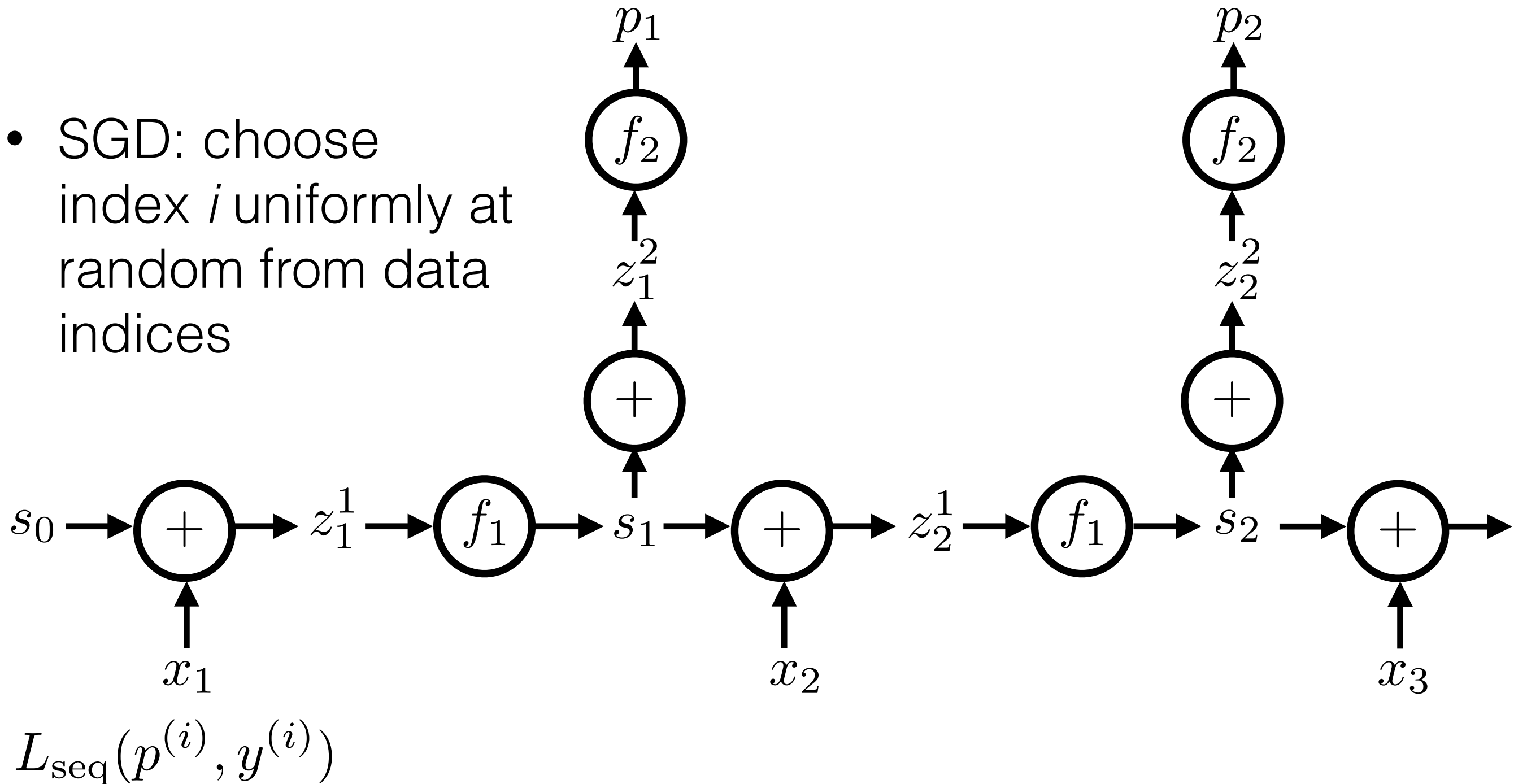
RNNs: a taste of backpropagation

- SGD: choose index i uniformly at random from data indices



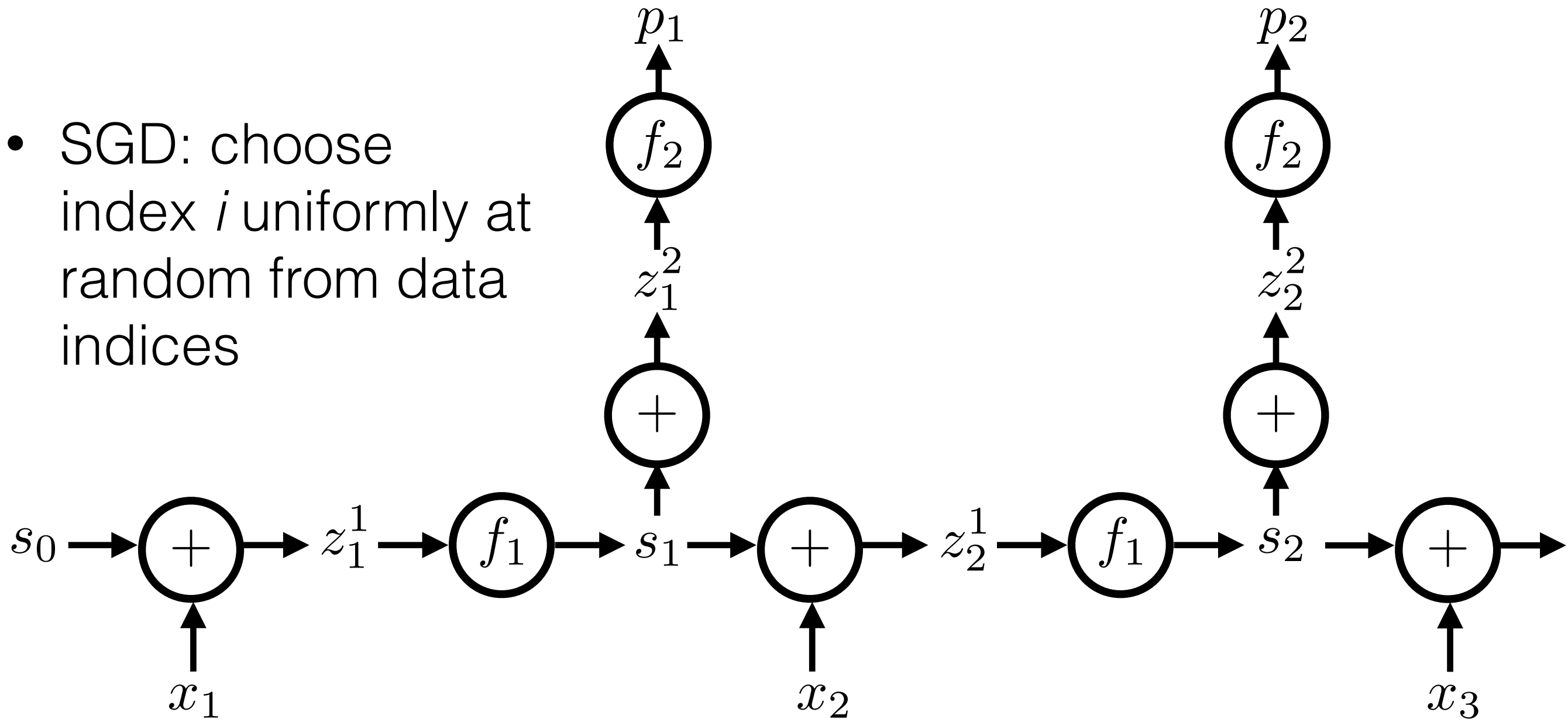
RNNs: a taste of backpropagation

- SGD: choose index i uniformly at random from data indices



RNNs: a taste of backpropagation

- SGD: choose index i uniformly at random from data indices

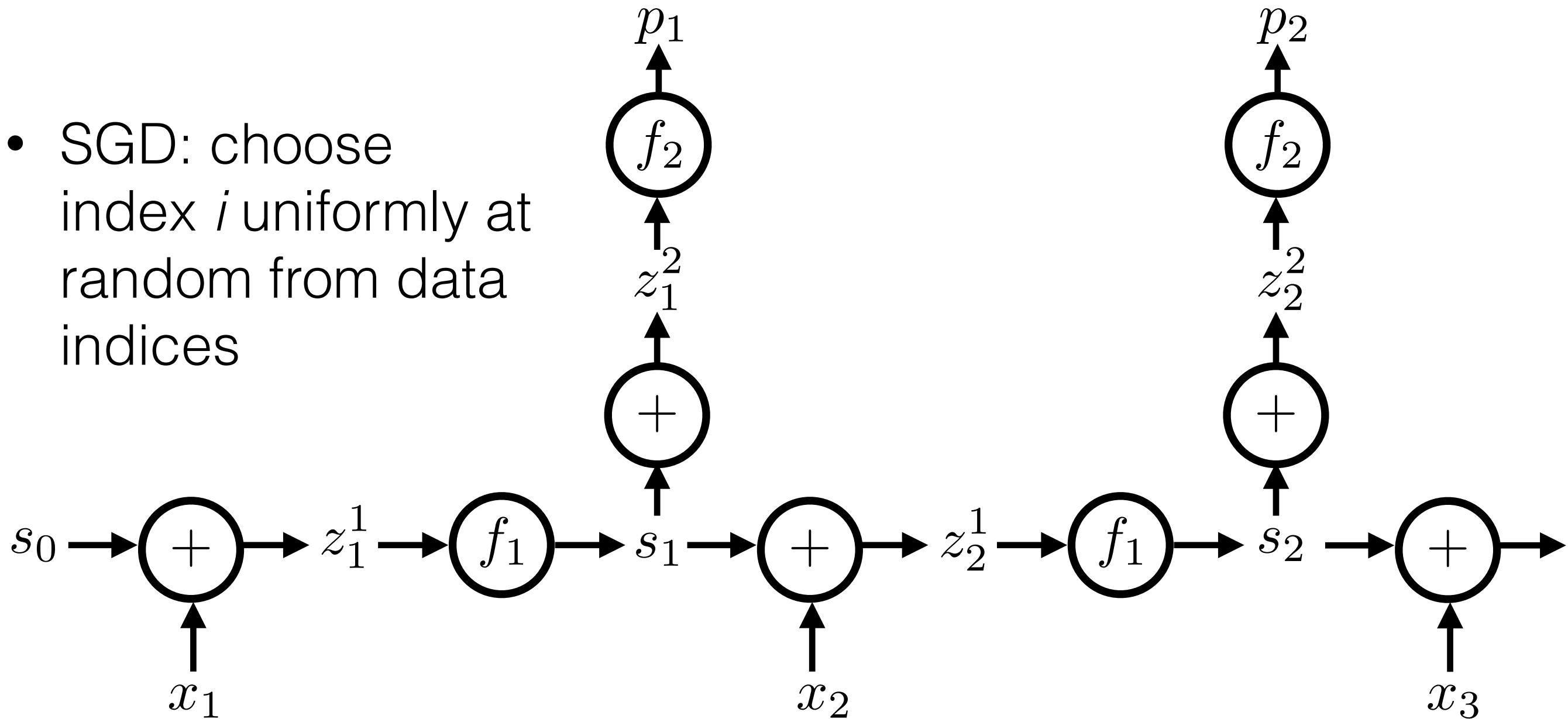


$$L_{\text{seq}}(p^{(i)}, y^{(i)})$$

$$\frac{dL_{\text{seq}}(p^{(i)}, y^{(i)})}{dW^{sx}}$$

RNNs: a taste of backpropagation

- SGD: choose index i uniformly at random from data indices

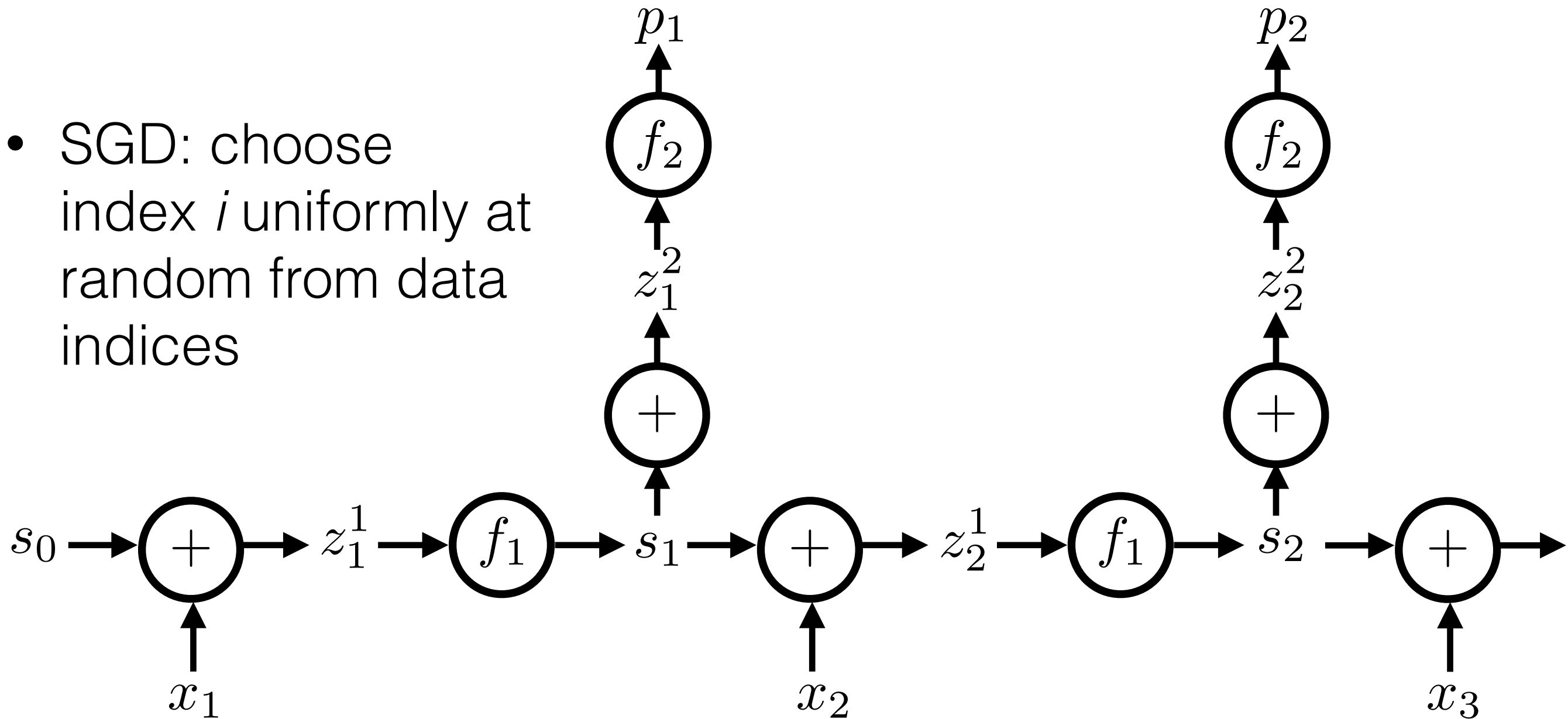


$$L_{\text{seq}}(p^{(i)}, y^{(i)})$$

$$\frac{dL_{\text{seq}}(p^{(i)}, y^{(i)})}{dW^{sx}}$$

RNNs: a taste of backpropagation

- SGD: choose index i uniformly at random from data indices

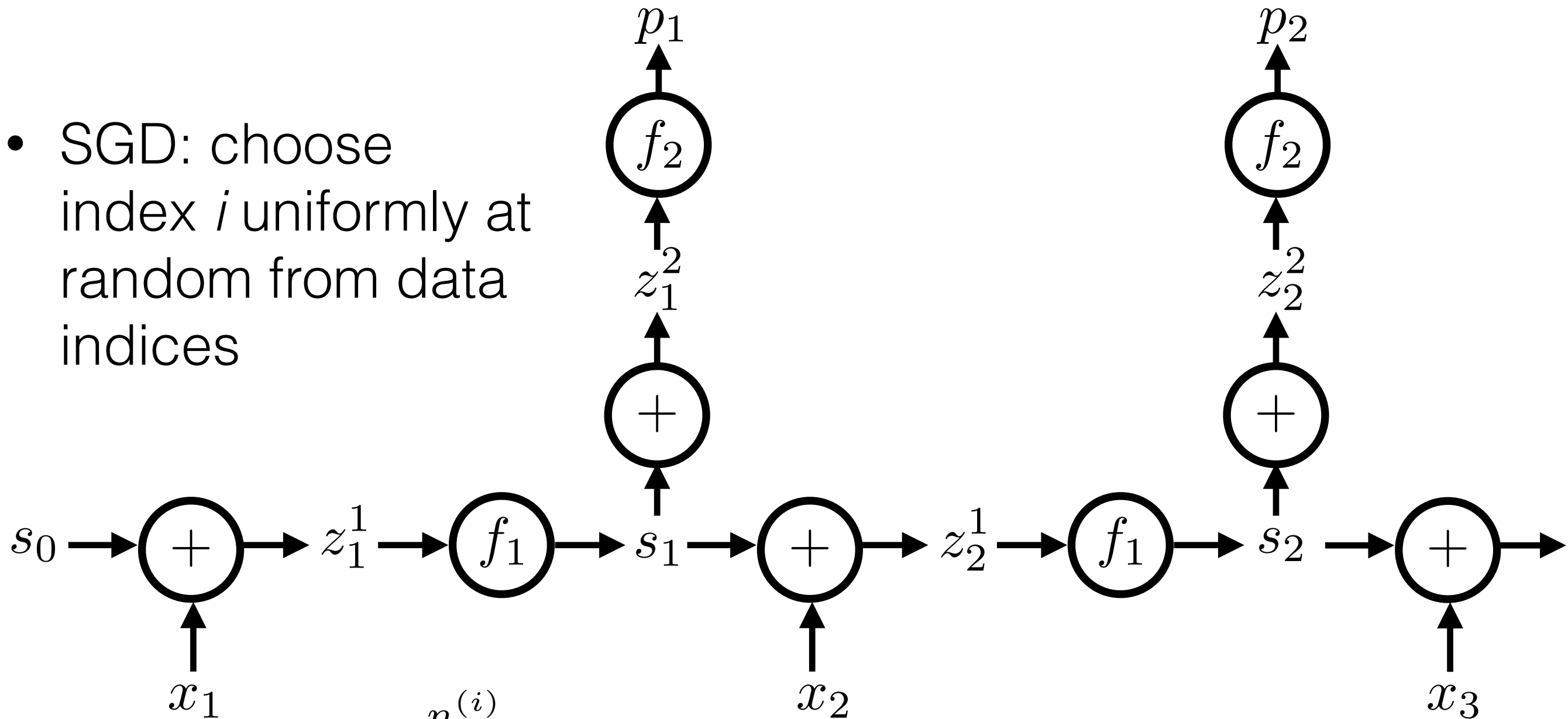


$$L_{\text{seq}}(p^{(i)}, y^{(i)})$$

$$\frac{dL_{\text{seq}}(p^{(i)}, y^{(i)})}{dW^{sx}}$$

RNNs: a taste of backpropagation

- SGD: choose index i uniformly at random from data indices

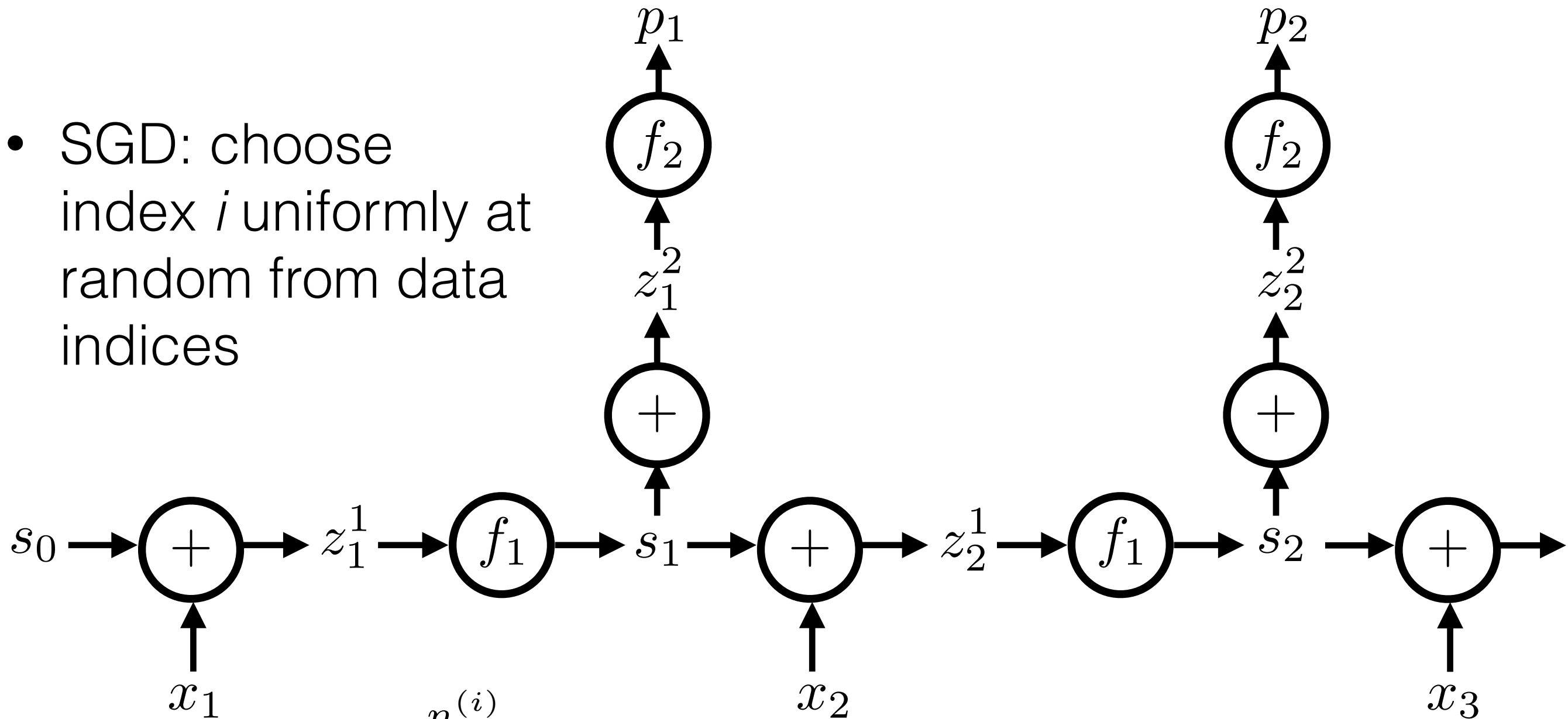


$$L_{\text{seq}}(p^{(i)}, y^{(i)}) = \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)})$$

$$\frac{dL_{\text{seq}}(p^{(i)}, y^{(i)})}{dW^{sx}}$$

RNNs: a taste of backpropagation

- SGD: choose index i uniformly at random from data indices

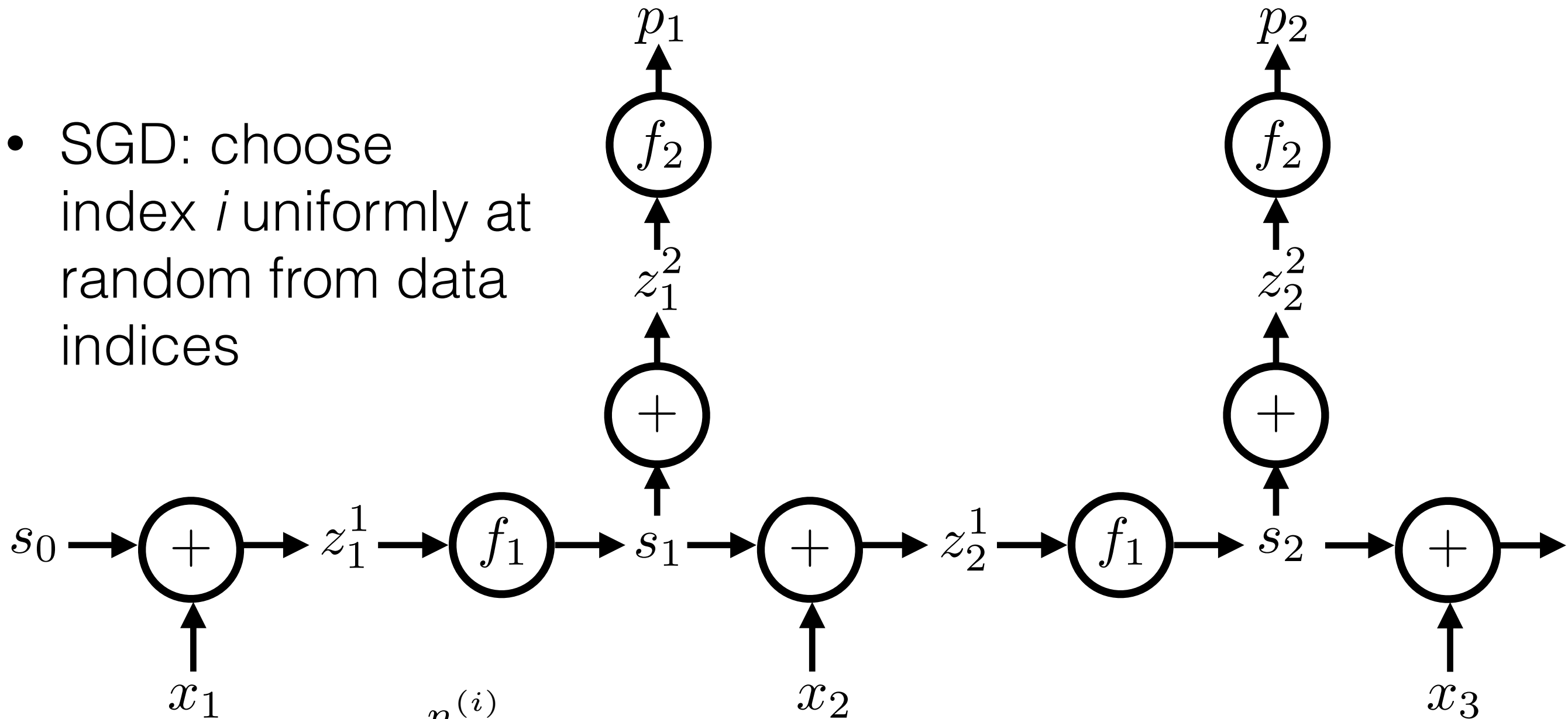


$$L_{\text{seq}}(p^{(i)}, y^{(i)}) = \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)})$$

$$\frac{dL_{\text{seq}}(p^{(i)}, y^{(i)})}{dW^{sx}} = \sum_{t=1}^{n^{(i)}} \frac{dL_{\text{elt}}(p_t^{(i)}, y_t^{(i)})}{dW^{sx}}$$

RNNs: a taste of backpropagation

- SGD: choose index i uniformly at random from data indices



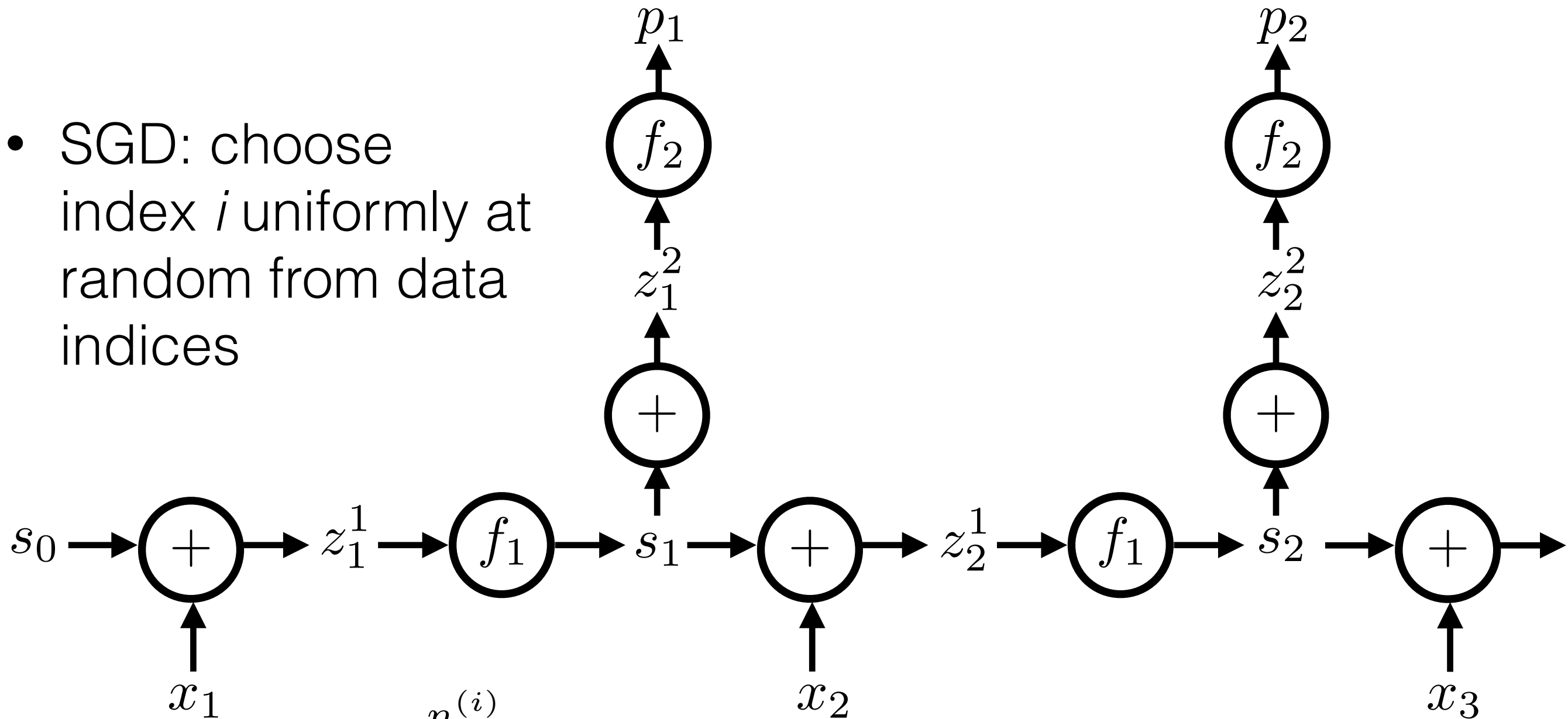
$$L_{\text{seq}}(p^{(i)}, y^{(i)}) = \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)})$$

$$L_t := L_{\text{elt}}(p_t^{(i)}, y_t^{(i)})$$

$$\frac{dL_{\text{seq}}(p^{(i)}, y^{(i)})}{dW^{sx}} = \sum_{t=1}^{n^{(i)}} \frac{dL_{\text{elt}}(p_t^{(i)}, y_t^{(i)})}{dW^{sx}}$$

RNNs: a taste of backpropagation

- SGD: choose index i uniformly at random from data indices



$$L_{\text{seq}}(p^{(i)}, y^{(i)}) = \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)})$$

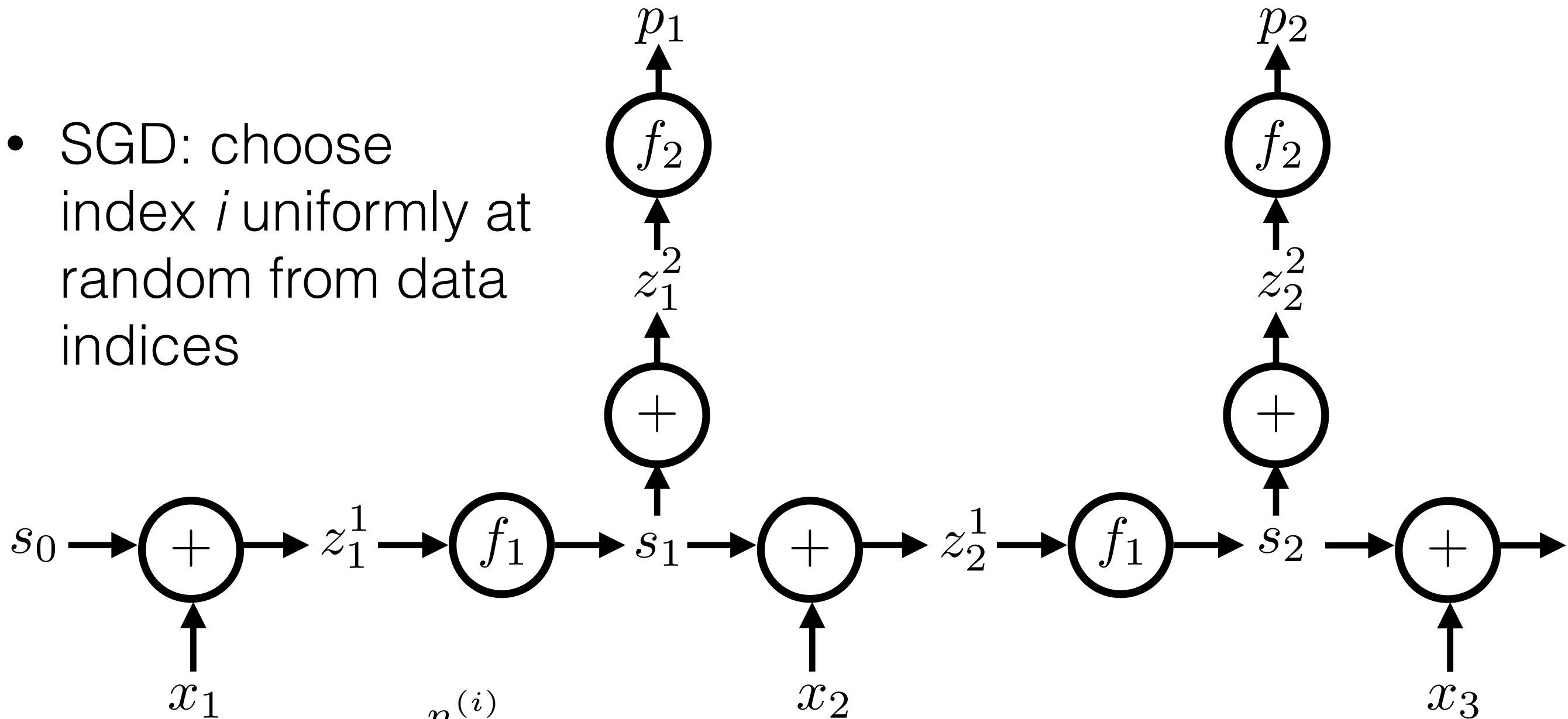
$$\frac{dL_{\text{seq}}(p^{(i)}, y^{(i)})}{dW^{sx}} = \sum_{t=1}^{n^{(i)}} \frac{dL_{\text{elt}}(p_t^{(i)}, y_t^{(i)})}{dW^{sx}}$$

$$L_t := L_{\text{elt}}(p_t^{(i)}, y_t^{(i)})$$

- Need: $\frac{dL_t}{dW^{sx}}$

RNNs: a taste of backpropagation

- SGD: choose index i uniformly at random from data indices



$$L_{\text{seq}}(p^{(i)}, y^{(i)}) = \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)})$$

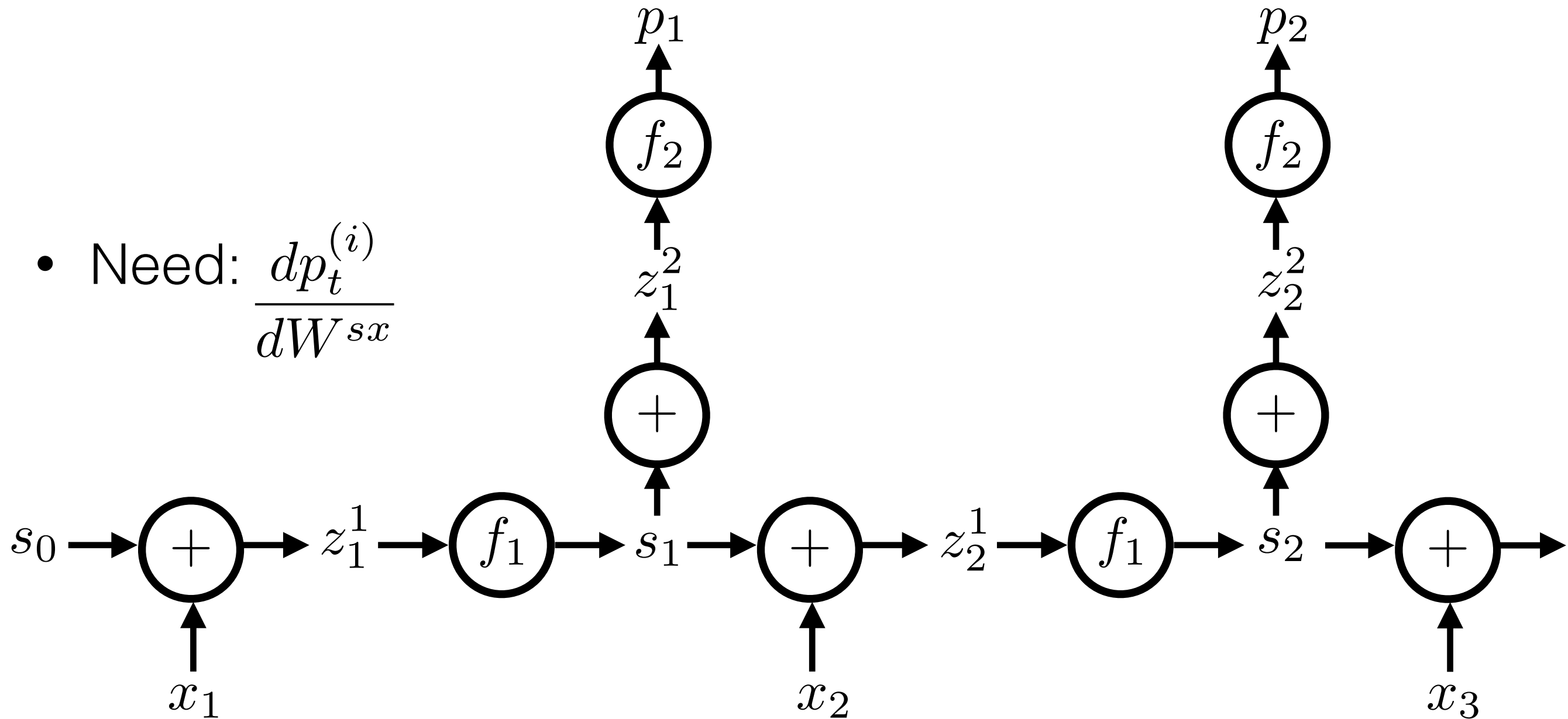
$$\frac{dL_{\text{seq}}(p^{(i)}, y^{(i)})}{dW^{sx}} = \sum_{t=1}^{n^{(i)}} \frac{dL_{\text{elt}}(p_t^{(i)}, y_t^{(i)})}{dW^{sx}}$$

$$L_t := L_{\text{elt}}(p_t^{(i)}, y_t^{(i)})$$

- Need: $\frac{dp_t^{(i)}}{dW^{sx}}$

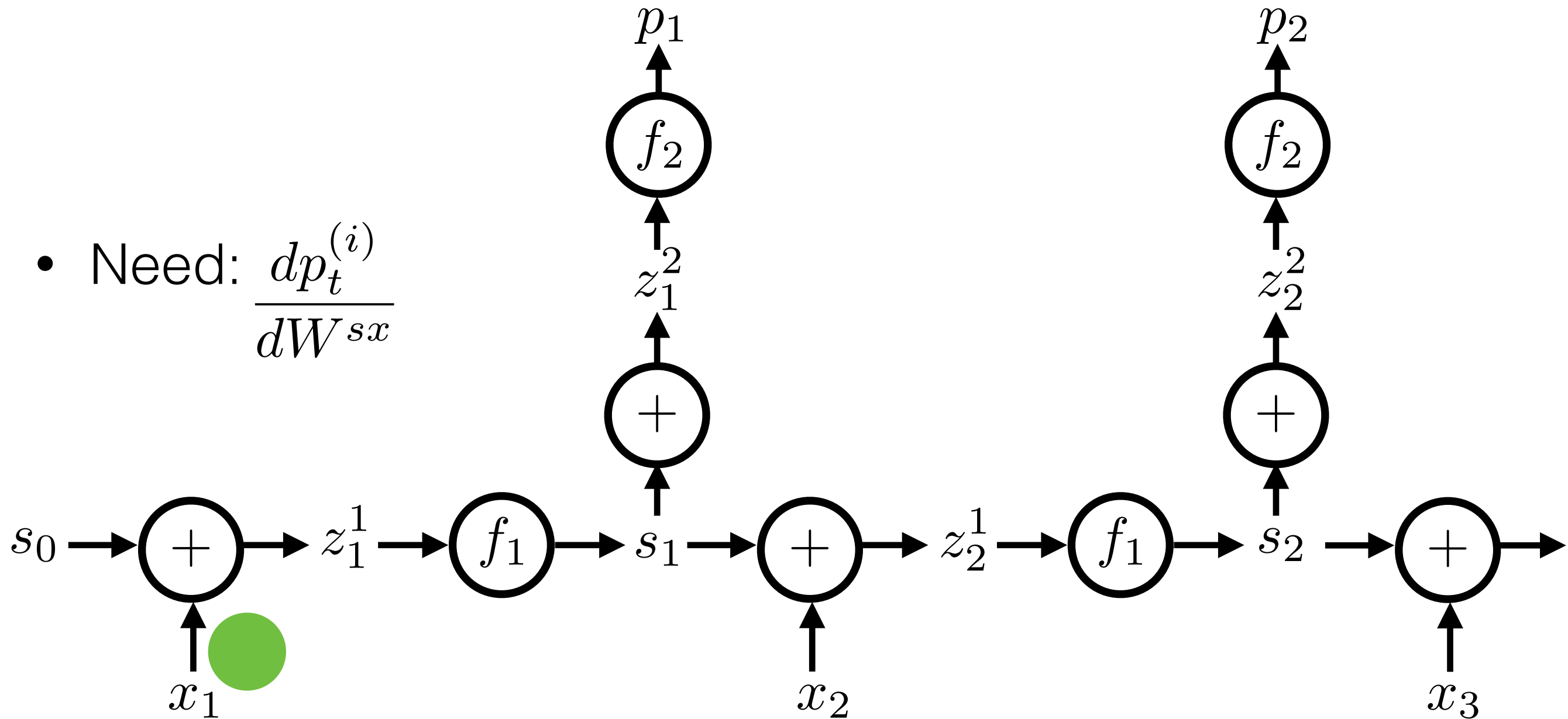
RNNs: a taste of backpropagation

- Need: $\frac{dp_t^{(i)}}{dW^{sx}}$



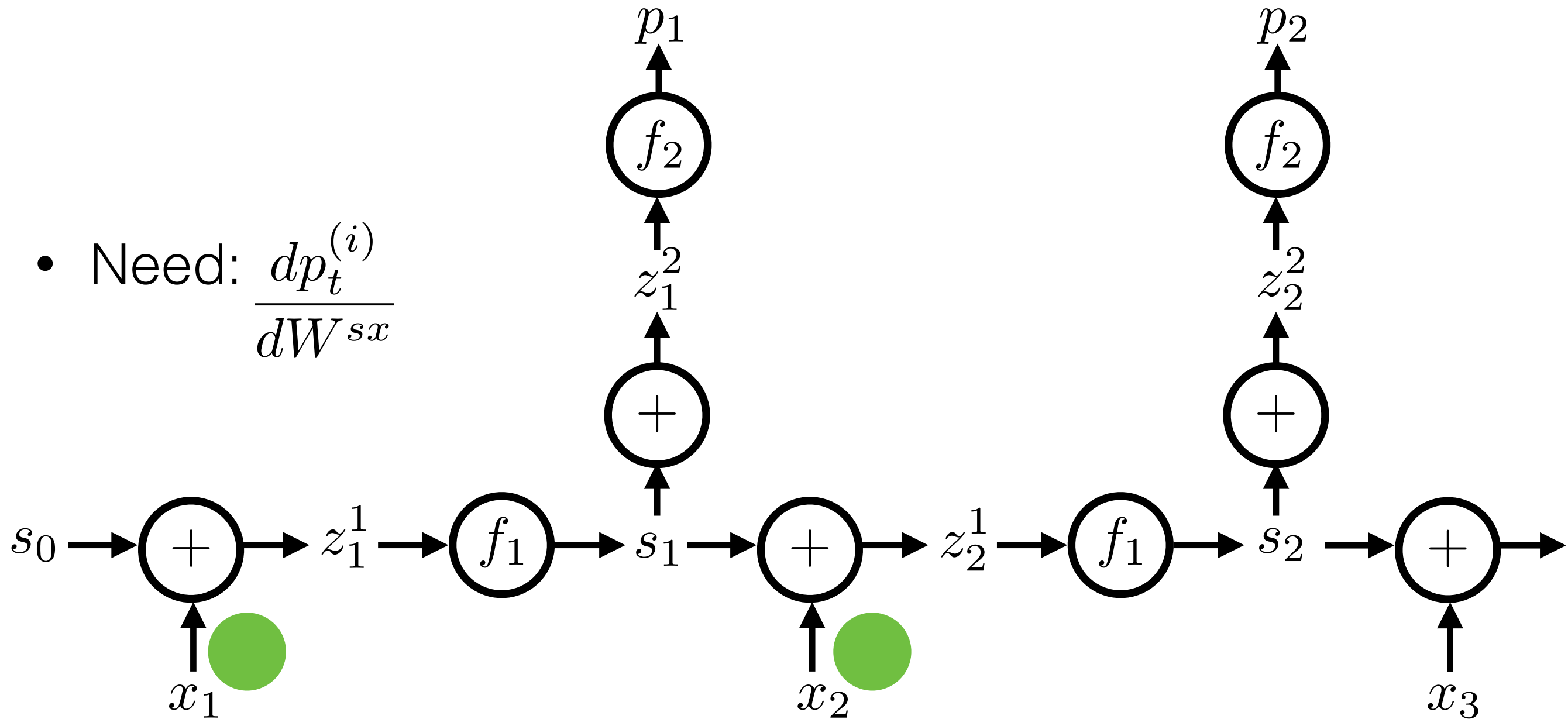
RNNs: a taste of backpropagation

- Need: $\frac{dp_t^{(i)}}{dW^{sx}}$



RNNs: a taste of backpropagation

- Need: $\frac{dp_t^{(i)}}{dW^{sx}}$



RNNs: a taste of backpropagation

- Need: $\frac{dp_t^{(i)}}{dW^{sx}}$

