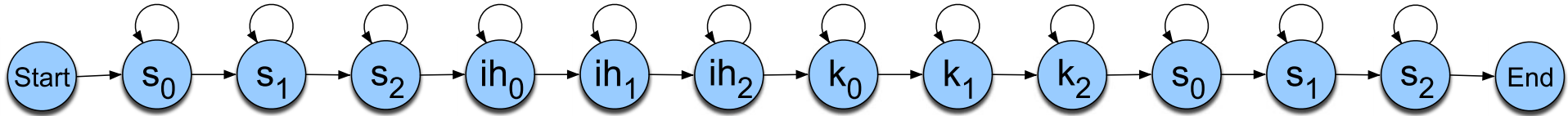


Распознавание речи

Современные подходы

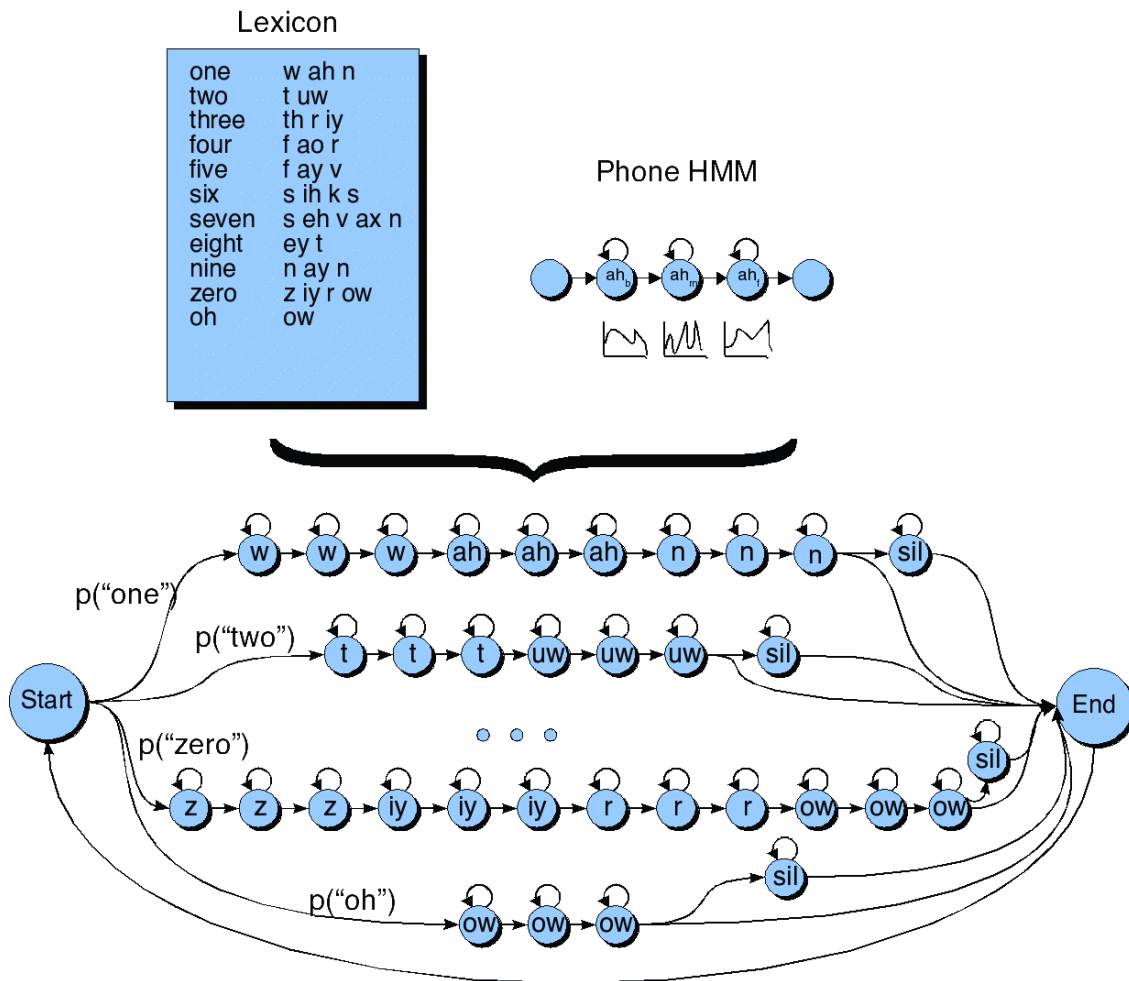
Декодирование для чисел

- НММ для одного слова



- $Q = q_1, \dots, q_N = N$ – состояний модели – это наши субфонемы
- $A = a_{11}, a_{12}, a_{1n}, \dots, a_{nn}$, - матрица вероятностей перехода. a_{ij} - вероятность перехода между субфонемами или переход сам на себя. Для каждого слова создается граф состояний НММ.
- $V = b_j(o_t)$, – набор наблюдений. Каждый элемент элемент выражает вероятность пронаблюдать вектор признаков, который генерируется состоянием j .

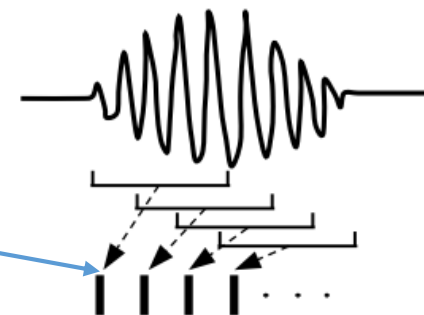
Декодирование для изолированных слов



- Матрицы A и B получаем в результате обучения модели алгоритмом “Вперед-назад”.

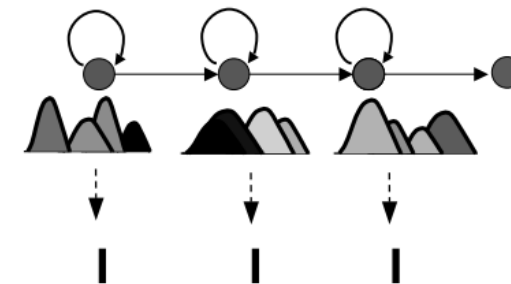
Декодирование для изолированных слов

- $\hat{W} = \operatorname{argmax}_{W \in \mathcal{L}} P(O|W)P(W)$



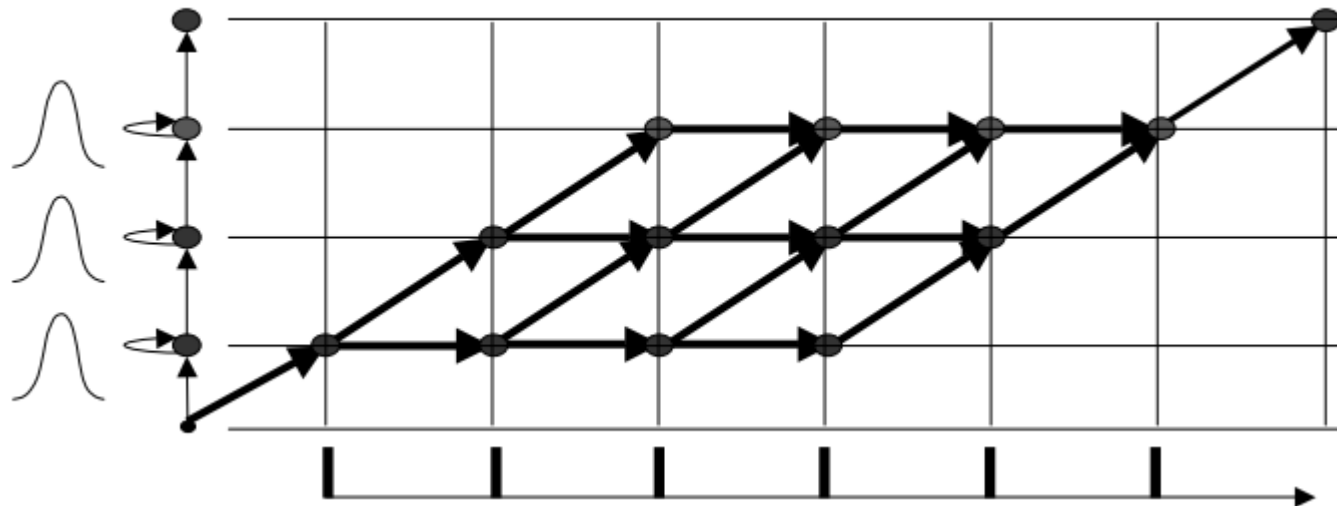
- $P(W)$ - априорная вероятность слова
- $P(O|W)$ - вероятность наблюдения признаков сгенерированных словом W
- Мы представляем каждое слово с помощью НММ

- НММ - это графическая форма функции плотности вероятности для изменяющихся во времени данных



$P(O|W)$ - для изолированного слова

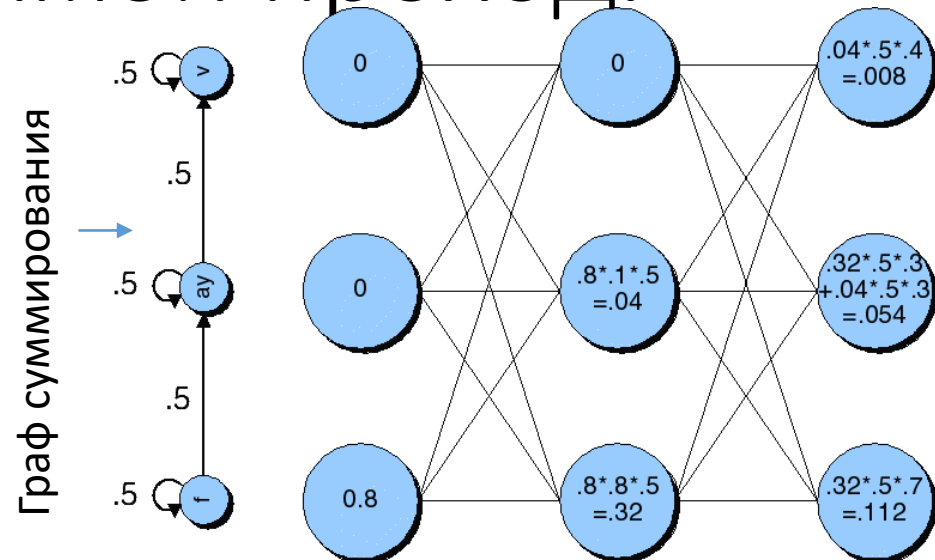
- *Реальная последовательность состояний, которое генерирует O неизвестно*
 - $P(O|W)$ - должно рассматривать все возможные состояния
 - $P(O|W) = \sum_{q \in \{\text{все состояния}\}} P(O, q|W)$



Вероятность
последовательности это
сумма вероятностей по
всем последовательностям

$P(O|W)$ - для слова "five". Сумма по всем ВОЗМОЖНЫМ путям. Прямой проход.

- f ay ay ay ay v v v v
- f f ay ay ay ay v v v
- f f f ay ay ay ay v
- f f ay ay ay ay ay ay v
- f f ay ay ay ay ay ay ay ay v
- f f ay v v v v v v v



V	0	0	0.008	0.0093	0.0114	0.00703	0.00345	0.00306	0.00206	0.00117
AY	0	0.04	0.054	0.0664	0.0355	0.016	0.00676	0.00208	0.000532	0.000109
F	0.8	0.32	0.112	0.0224	0.00448	0.000896	0.000179	4.48e-05	1.12e-05	2.8e-06
Time	1	2	3	4	5	6	7	8	9	10
B	<i>f</i>	0.8	<i>f</i> 0.8	<i>f</i> 0.7	<i>f</i> 0.4	<i>f</i> 0.4	<i>f</i> 0.4	<i>f</i> 0.5	<i>f</i> 0.5	<i>f</i> 0.5
	<i>ay</i>	0.1	<i>ay</i> 0.1	<i>ay</i> 0.3	<i>ay</i> 0.8	<i>ay</i> 0.8	<i>ay</i> 0.8	<i>ay</i> 0.6	<i>ay</i> 0.5	<i>ay</i> 0.4
	<i>v</i>	0.6	<i>v</i> 0.6	<i>v</i> 0.4	<i>v</i> 0.3	<i>v</i> 0.3	<i>v</i> 0.3	<i>v</i> 0.6	<i>v</i> 0.8	<i>v</i> 0.9
	<i>p</i>	0.4	<i>p</i> 0.4	<i>p</i> 0.2	<i>p</i> 0.1	<i>p</i> 0.1	<i>p</i> 0.1	<i>p</i> 0.1	<i>p</i> 0.3	<i>p</i> 0.3
	<i>iy</i>	0.1	<i>iy</i> 0.1	<i>iy</i> 0.3	<i>iy</i> 0.6	<i>iy</i> 0.6	<i>iy</i> 0.6	<i>iy</i> 0.6	<i>iy</i> 0.5	<i>iy</i> 0.5

наблюдения

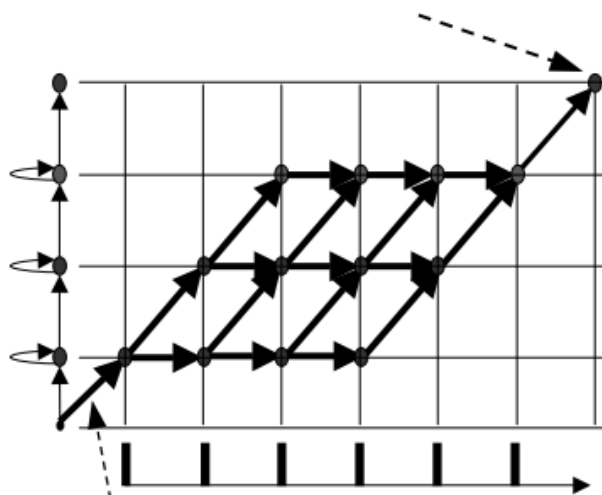
Классификация изолированных слов

Выбираем:
 $\operatorname{argmax}_{W \in \mathcal{L}} P(O|W)P(W)$

HMM для Odd



$$P(\text{Odd})P(O|\text{Odd})$$

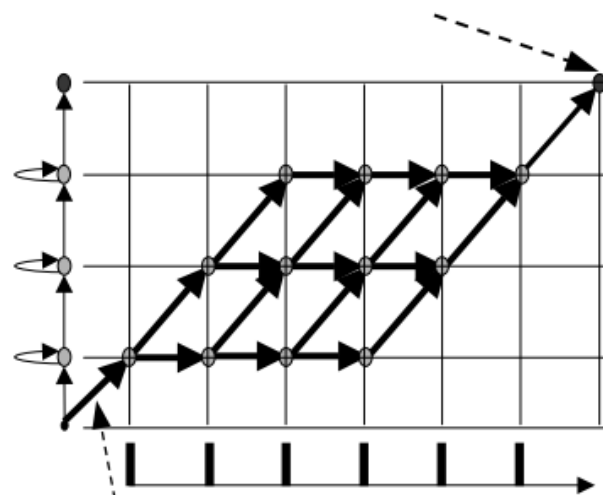


$$P(\text{Odd})$$

HMM для Even

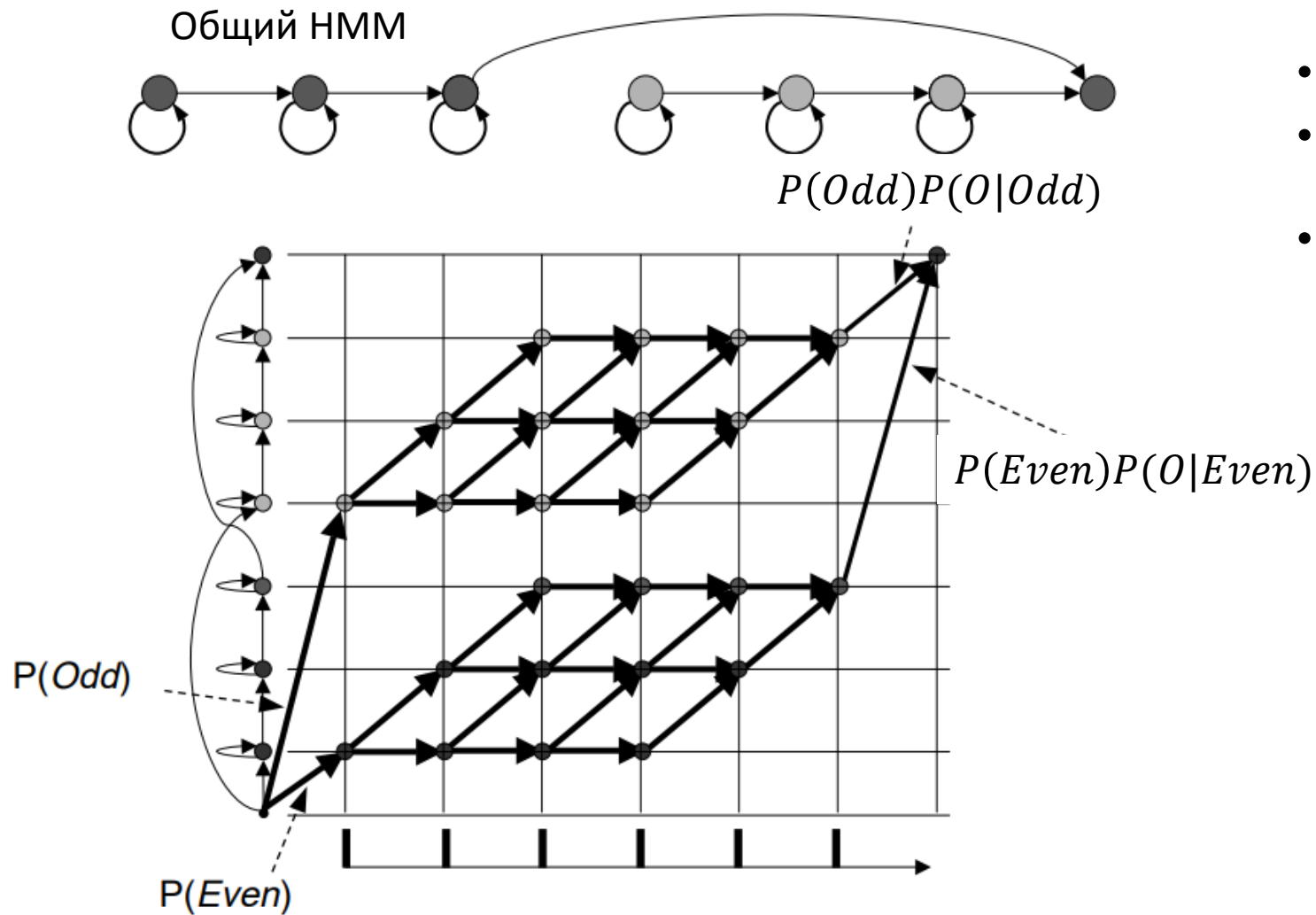


$$P(\text{Even})P(O|\text{Even})$$



$$P(\text{Even})$$

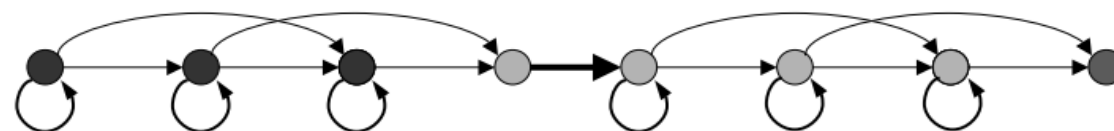
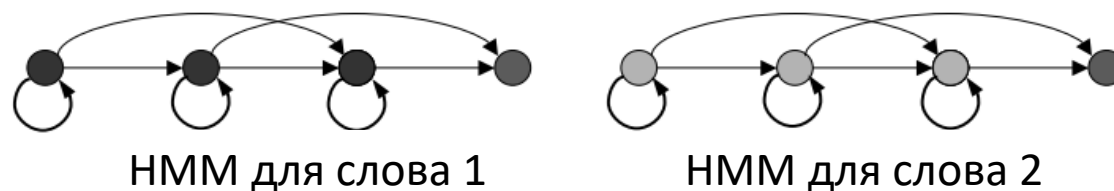
Классификация между Odd и Even



- Используем общий НММ
- Для одной фонемы используем одно состояние
- Для декодирования используем Viterbi.
- $\max\{P(Odd)P(O|Odd), P(Even)P(O|Even)\}$

HMM для последовательности слов

- Складываем HMM для каждого слова в последовательности

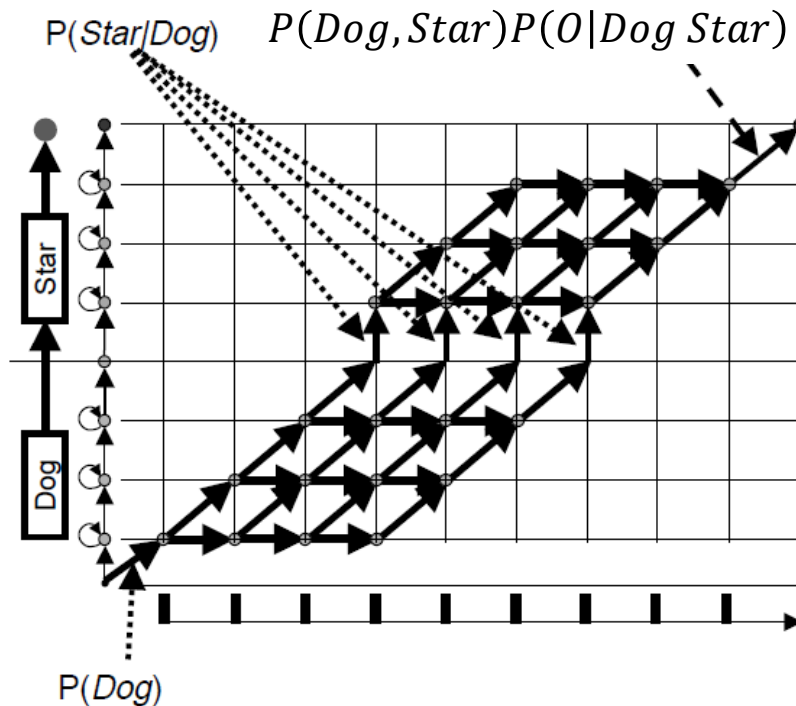
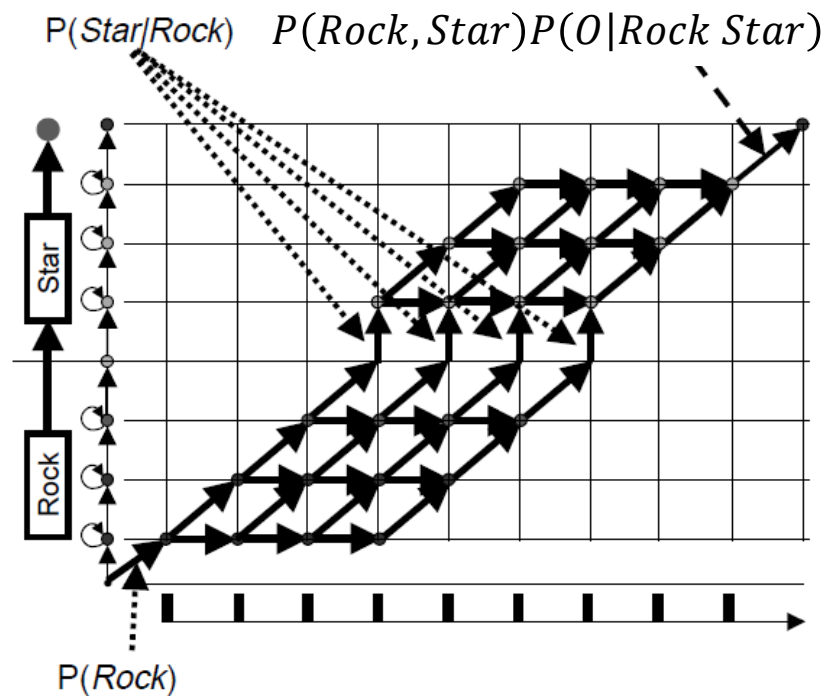


Комбинированный HMM для последовательности слово 1 – слово 2

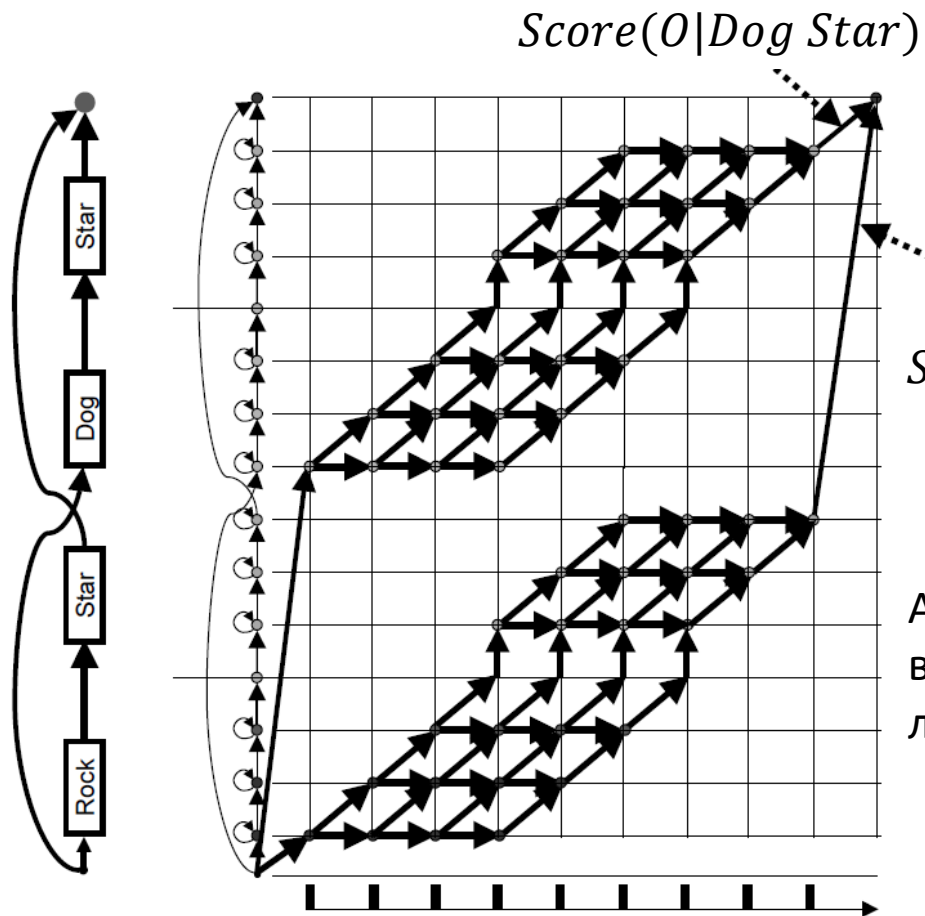
- HMM для слов - сама по себе конкатенация HMM для фонем
 - фонем намного меньше, чем слов, и встречаются чаще в данных обучения
 - слова, которые никогда не видели в данных обучения, могут быть построены из фонемных HMM

Классификация двух последовательностей

- Классифицируем “Rock Star” vs “Dog Star”
 - Прямой проход и выбираем максимум:
 - $\text{argmax} \{P(\text{Rock}, \text{Star})P(O|\text{Rock Star}), P(\text{Dog}, \text{Star})P(O|\text{Dog Star})\}$



Декодирование последовательности слов

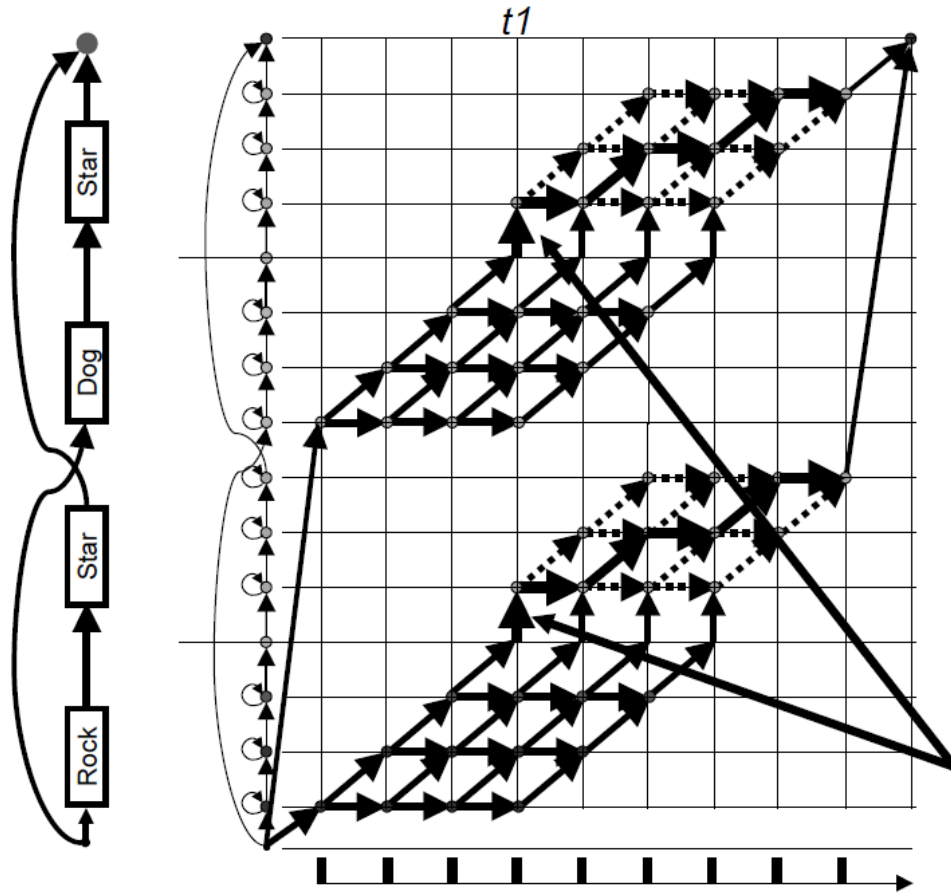


Используем Viterbi для получения сора
каждого пути
 $Viterbi = \max(\max(dogstar), \max(rockstar))$

$Score(O|Rock Star)$

Аппроксимируем общую
вероятность выбором
лучшего пути

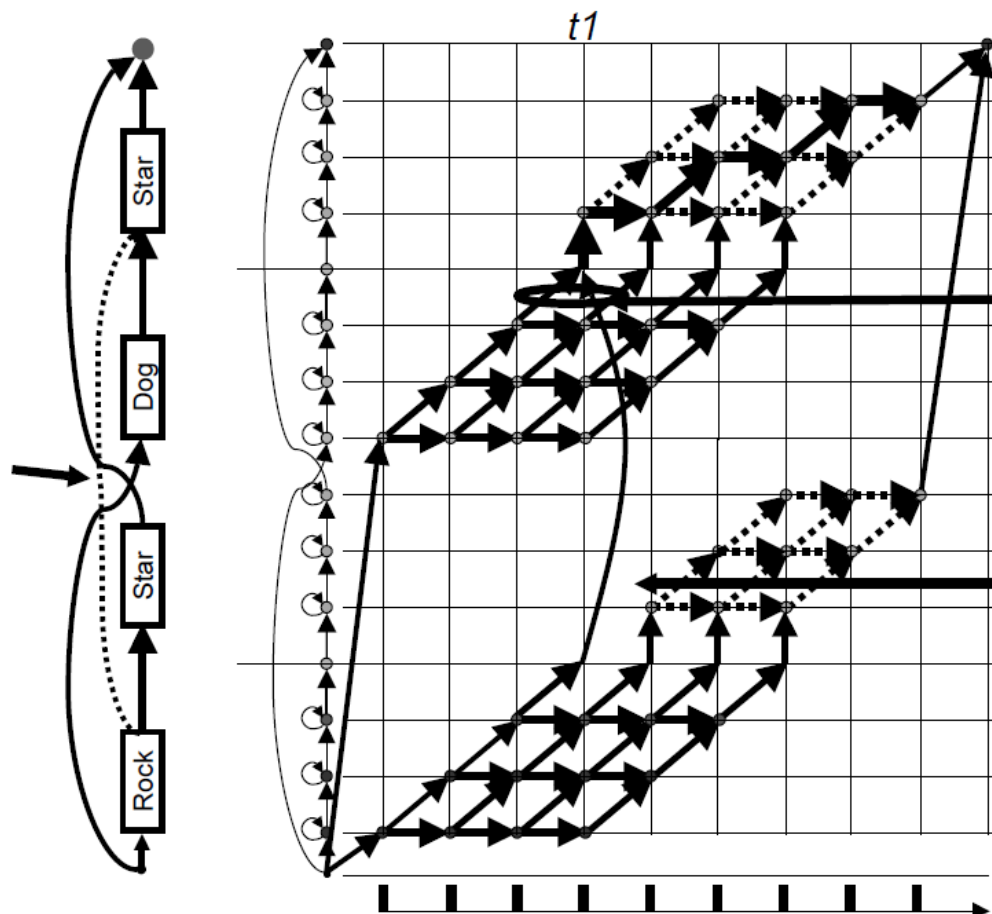
Декодирование последовательностей



Для финальной точки наилучший путь через "STAR" одинаковый для обеих решеток

Мы можем выбирать между "Dog" и "Rock" прямо здесь, потому что фичи этих путей идентичны.

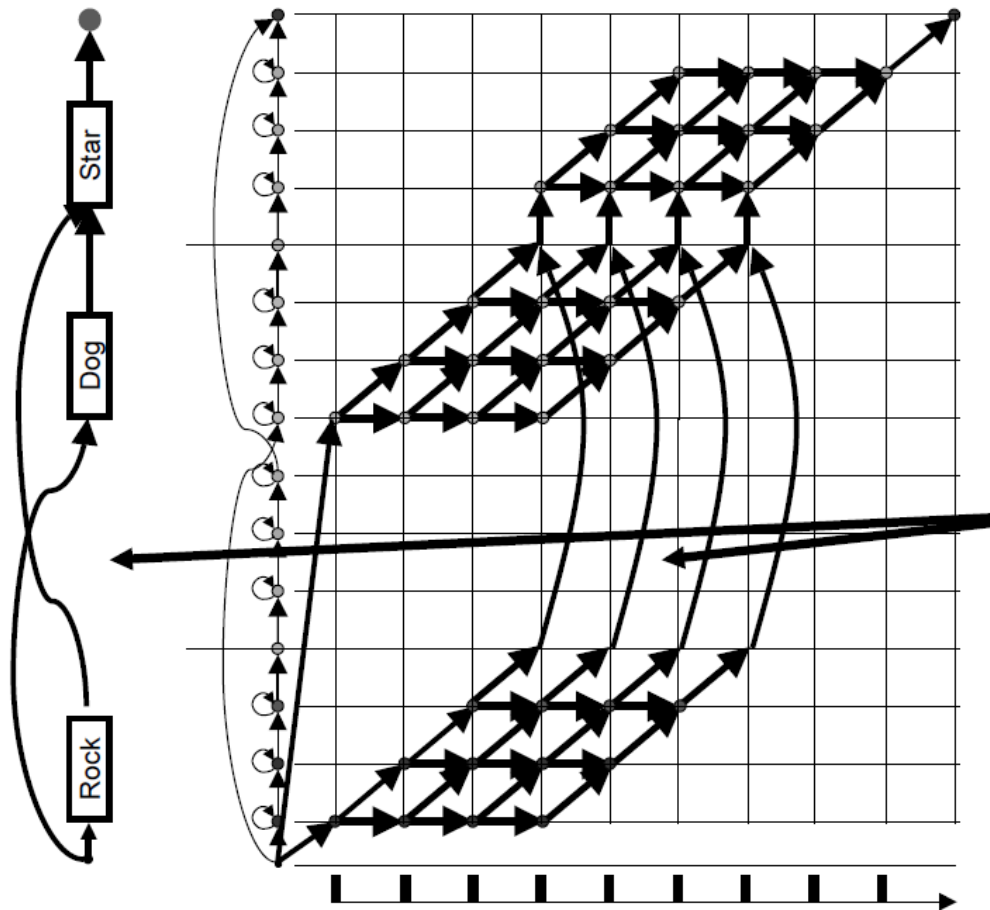
Декодирование последовательностей



Мы выбираем наилучший
скор двух входящих ребер
тут

Эту часть решетки удаляем

Декодирование последовательности



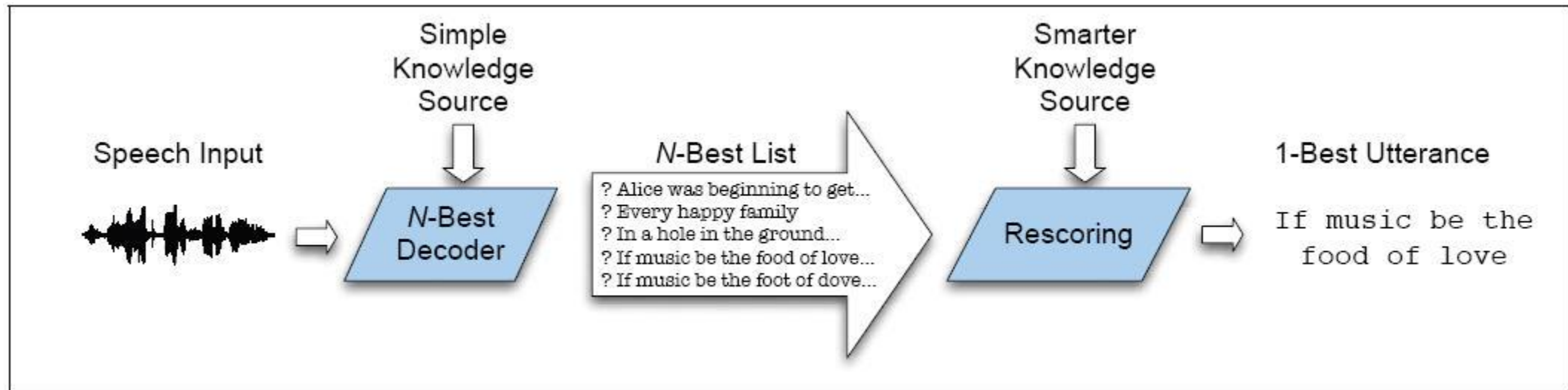
Мы теперь выполняем декодирование Витерби на этом свернутом графе.

Эту копию решетки для "STAR" полностью удаляем

Проблемы с Витерби

- Трудно интегрировать сложные источники знаний
 - Триграммные грамматики
 - Парсерное решение или нейронную сеть для языковой модели
 - длинные зависимости нарушают предположения динамического программирования
 - Знания, которые идут не справа налево
 - Следующие слова могут помочь предсказать предыдущие слова
- Решения
 - Возвращать множество гипотез и затем ранжировать их по дополнительным знаниям
 - Использовать другой поисковый алгоритм

Многопроходный поиск

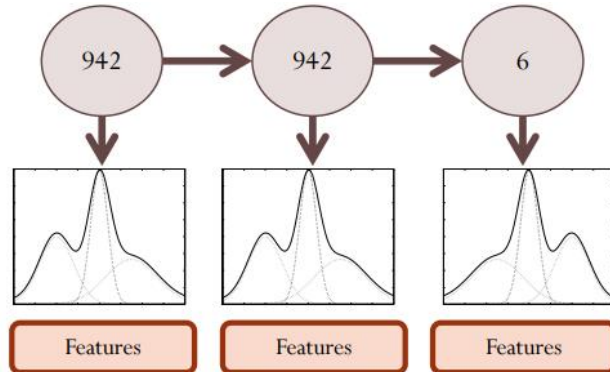


Расширенные алгоритмы поиска, разобрать самостоятельно.

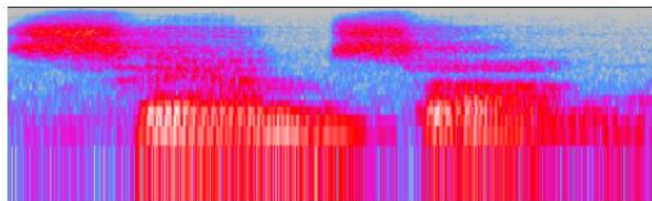
Нейросети для ASR

- Гибридная модель

Hidden Markov Model (HMM):



Акустическая модель:



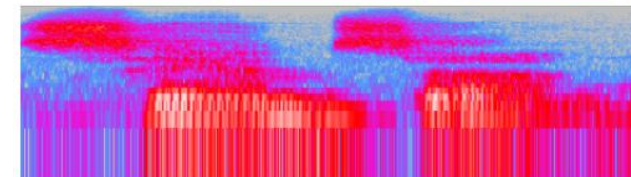
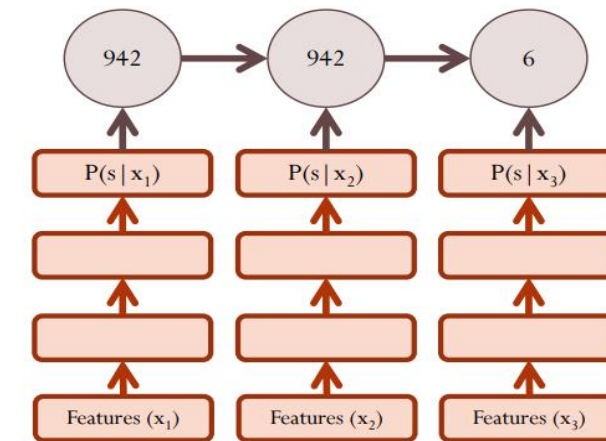
Аудио вход:

GMM Акустическая модель

Моделирует: $P(o|q)$

o : признаки

q : состояние HMM



Гибридная DNN акустическая модель

Моделирует: $P(q|o)$,

нам нужна $P(o|q) = P(q|o)P(o)/P(q)$

$P(o) = const$ – в процессе декодирования

$P(o|q) = P(q|o)/P(q)$

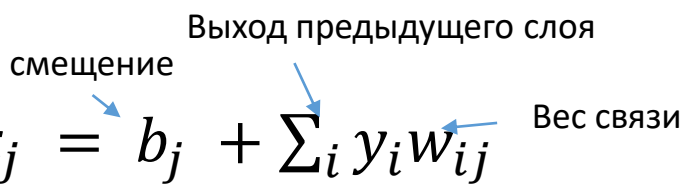
$P(q|o)$ - классификатор состояний

$P(q)$ – априорное распределение состояний

Целевая функция обучения DNN

- Обучение с учителем, минимизируем ошибку классификации
- Обычно – кросс энтропия
 - $Loss(x, y; W, b) = - \sum_{k=1}^K (y = k) \log f(x)_k$
- Ошибка считается по каждому окну.
- Метки классов получаем из HMM-GMM (принудительное выравнивание)

Обучение DNN

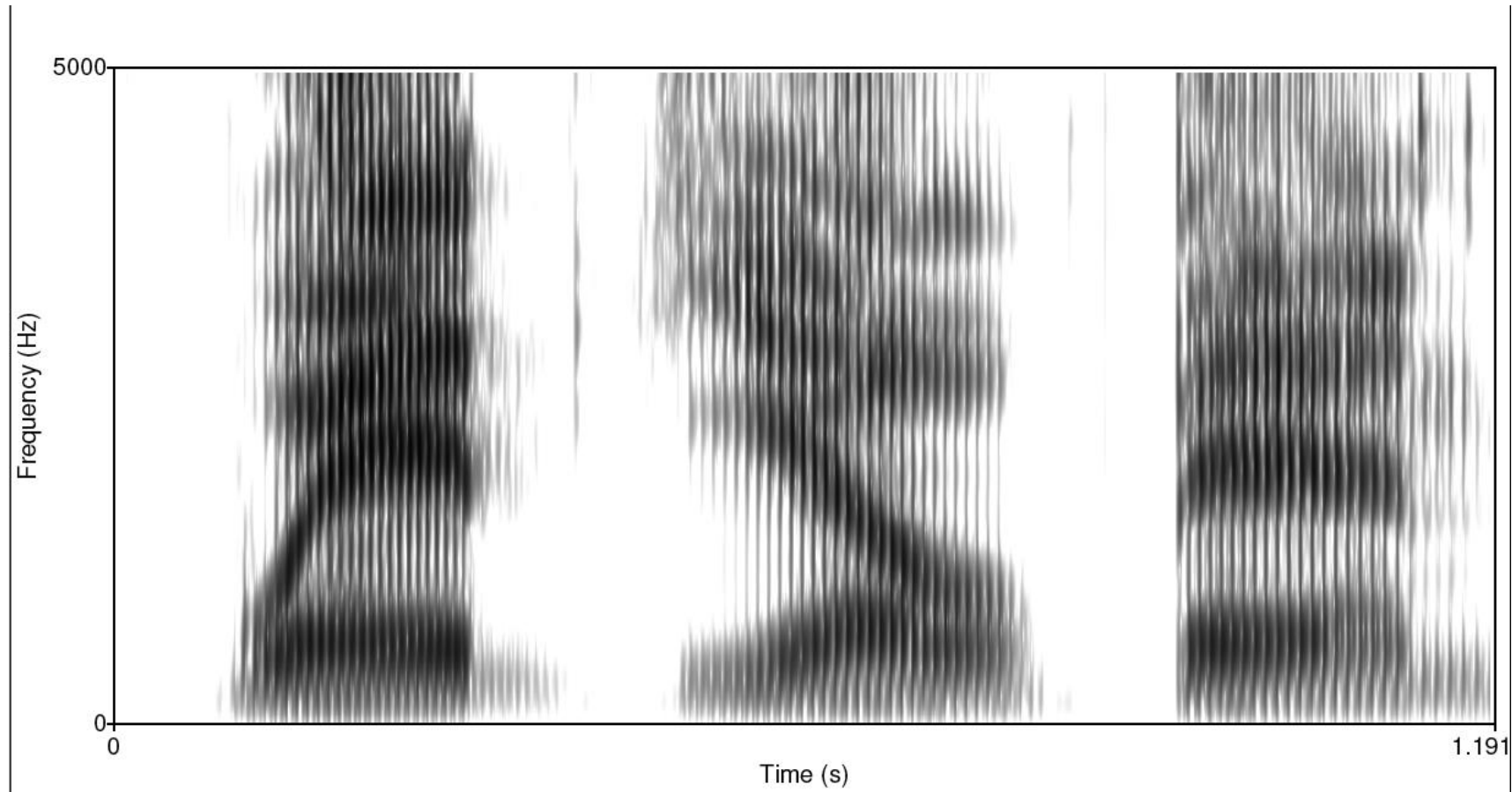
- Сеть состоит из более чем одного слоя между входом и выходом
 - Слой состоит из скрытых элементов
 - $j^{\text{й}}$ элемент слоя использует нелинейную функцию активации, например,
$$y_j = \text{logistic}(x_j) = \frac{1}{1 + e^{-x_j}}$$
 - Для маппирования его суммарного входа $x_j = b_j + \sum_i y_i w_{ij}$

- Для много классовой классификации выходы сети представляем вероятностями $p_j = \frac{\exp(x_j)}{\sum_k \exp(x_k)}$
- В качестве лоса используем кросс энтропию $L(x, y; W, b) = - \sum_{k=1}^K (y = k) \log f(x)_k$

Обучение DNN

- Две фазы обучения
 - Прямой проход
 - Обратный проход
- Прямой проход – подаем вход (фичи).
 - Рассчитываем активации на слоях и выходе сети
 - Рассчитываем полученную ошибку
- Обратный проход:
 - Считаем градиент ошибки по входу $\partial L / \partial x$
 - Используем цепное правило для распространения градиента по сети и обновления весов сети
 - $$\Delta w_{ij}(t) = \alpha \Delta w_{ij}(t - 1) - \epsilon \frac{\partial C}{\partial w_{ij}(t)}.$$

Неоднородность звучания фонем в разных словах

w eh d y eh l b eh n

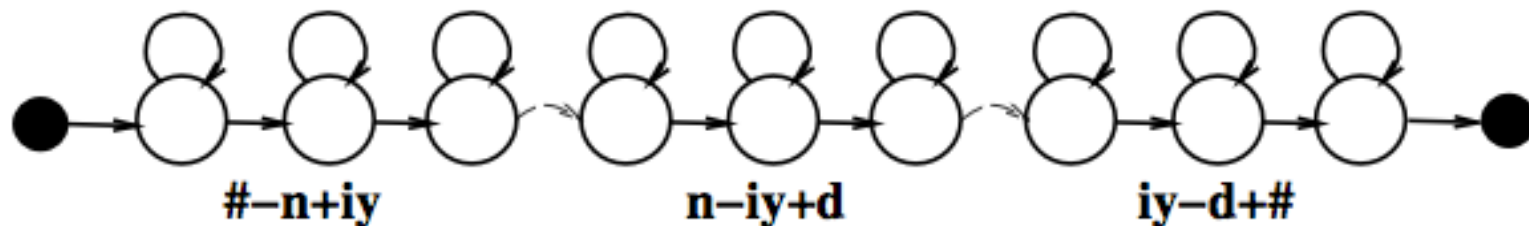


Моделирование фонетического контекста

- Самый сильный фактор, влияющий на фонетическую изменчивость это соседние фонемы
- Как моделировать это в НММ?
- Идея: иметь модели фонем, которые являются специфическими для контекста.
- Вместо контекстно-независимых (CI) фонем
- Мы будем использовать контекстно-зависимые (CD) фонемы

CD фонемы – трифоны (triphones)

- Каждый трифон инкапсулирует информацию о предыдущих и следующих фонемах
- Монофоны:
 - p, t, k
- Трифоны
 - iy-p+aa
 - a-b+c - значит “фонема a предшествует b за которой стоит c”



Нужно учитывать трифоны в модели

Моделирование границ слов

- Внутри словные контекстно зависимые (CD) модели

‘OUR LIST’:

SIL AA+R AA-R L+IH L-IH+S IH-S+T S-T

- Между словные контекстно зависимые (CD) модели

‘OUR LIST’:

SIL-AA+R AA-R+L R-L+IH L-IH+S IH-S+T S-T+SIL

- Учет между словных трифонов делает декодирование сложнее!

Последствия кросс-словных трифонов

- Если фонетический набор - 50 фонем, то возможное количество трифонов - $50^3 = 125000$.
- В реальной жизни не весь набор трифонов используется.
- В задаче WSI содержится примерно 20К слов
- Всего трифонов в обучающей сети 18500
- Нужна обобщающая модель!

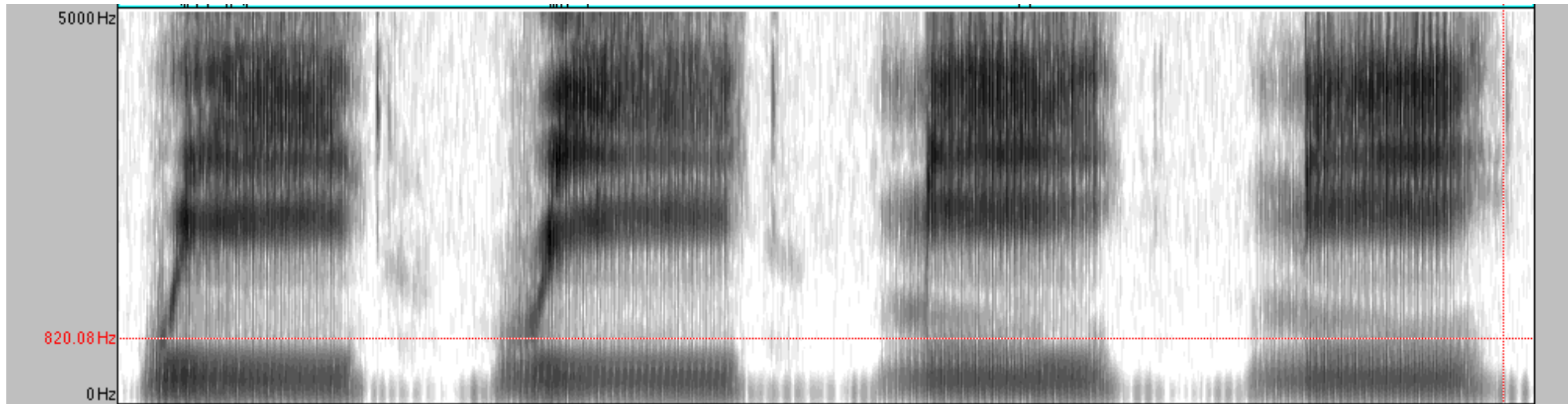
Моделирование фонетического контекста

w i y

r i y

m i y

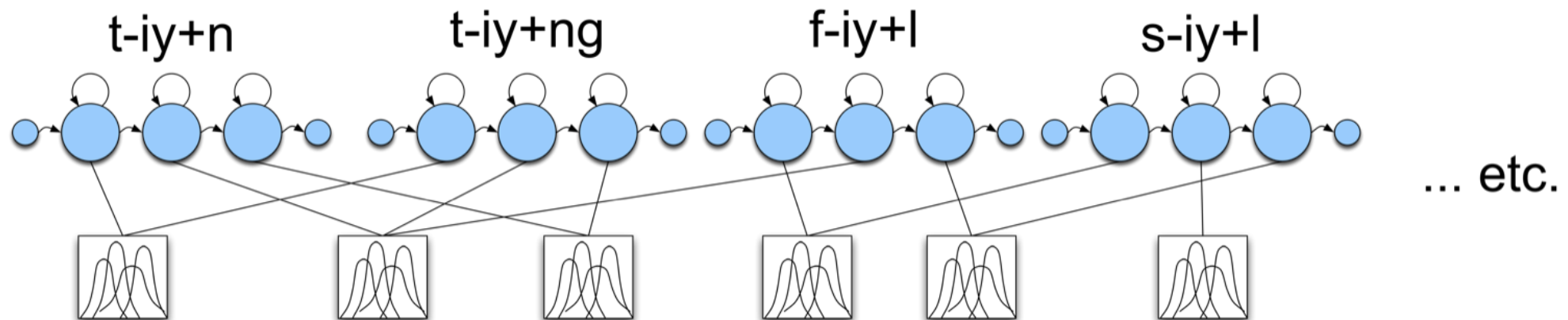
n i y



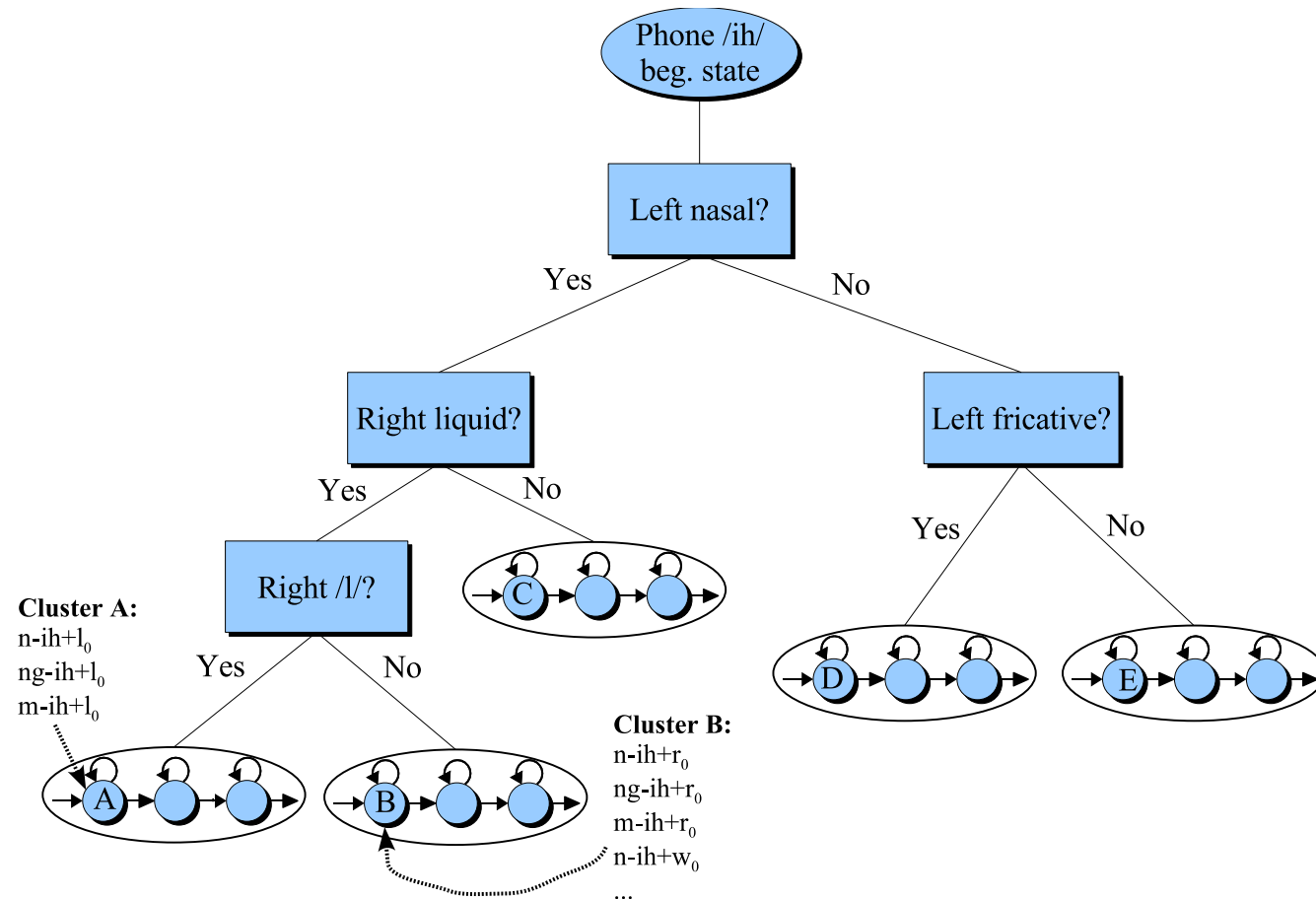
Некоторые фонемы в контексте выглядят одинаковыми

Связывание состояний *(Young et al)*

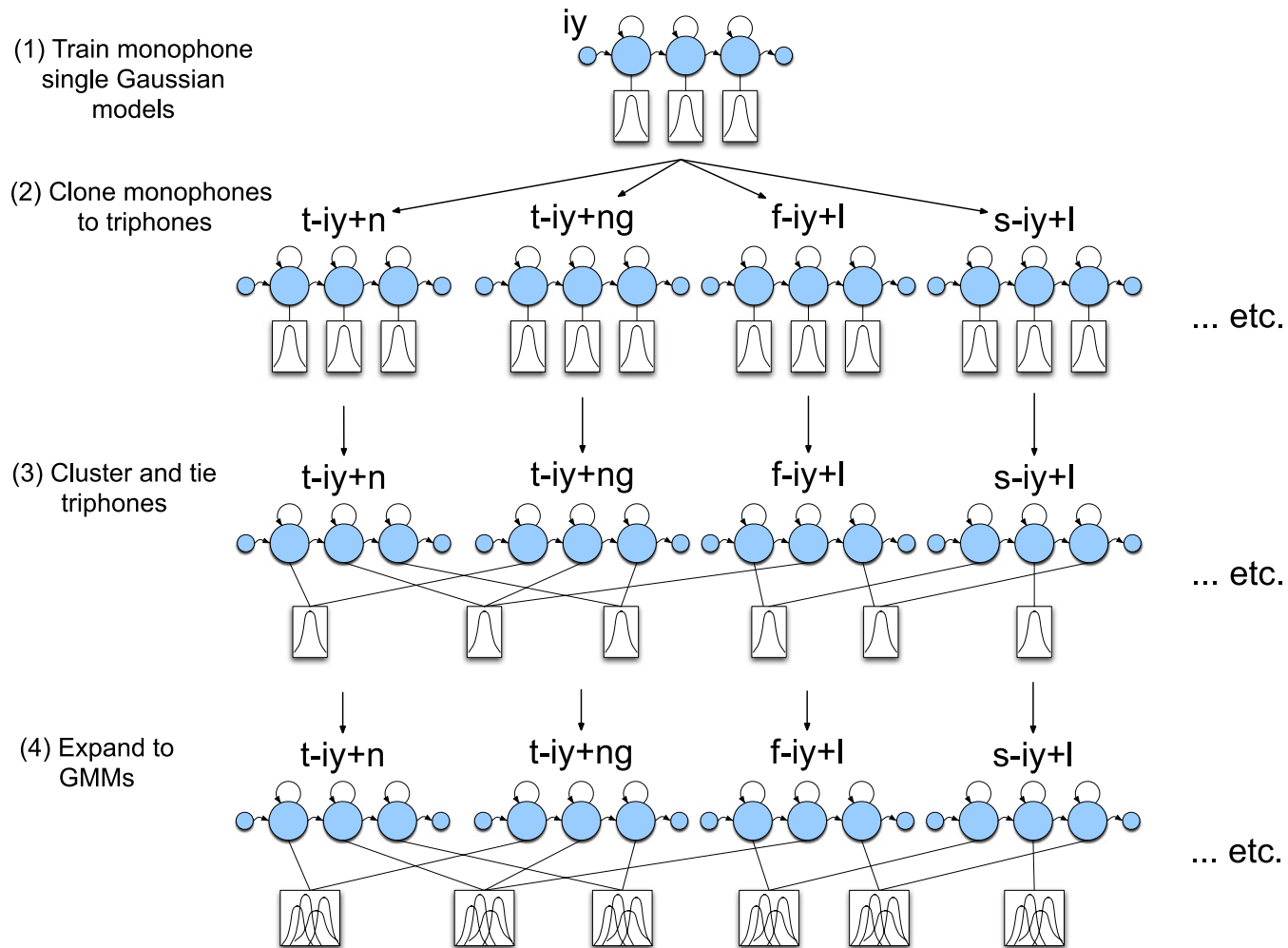
- Кластеризация трифонов на основе деревьев принятия решений
- Состояния, которые группируются вместе, будут разделять свои гауссианы
- Мы называем это «связывание состояний», поскольку эти состояния «связаны» с одним и тем же гауссовым распределением.



Дерево решений для СВЯЗЫВАНИЯ СОСТОЯНИЙ



Связывание состояний (State Tying: Young, Odell, Woodland 1994)



Нейросеть вместо GMM

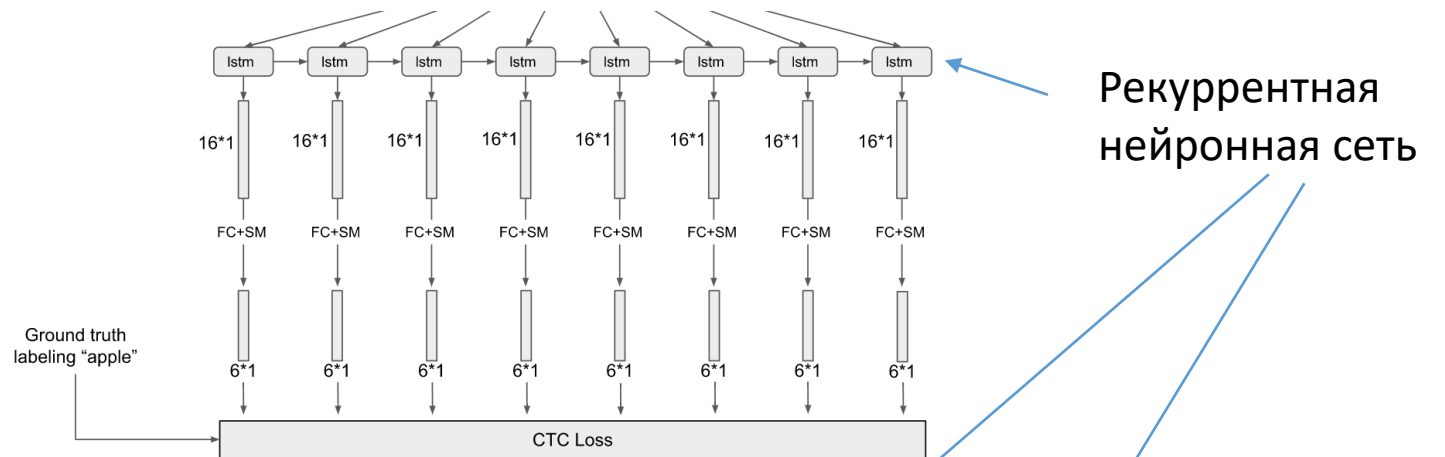
- Нейросеть лучший классификатор чем GMM, но нейросетью нельзя сделать выравнивание
- Используем нейросеть для классификации состояний и трифонов
 - Обучаем модель HMM + GMM + CD
 - На маркерах полученных из HMM+GMM +CD обучаем DNN
 - Еще раз делаем выравнивание при помощи HMM+DNN
 - F. Seide, G. Li, and D. Yu, "Conversational speech transcription using context-dependent deep neural networks."
 - Перетренируем DNN

Основные недостатки HMM + DNN

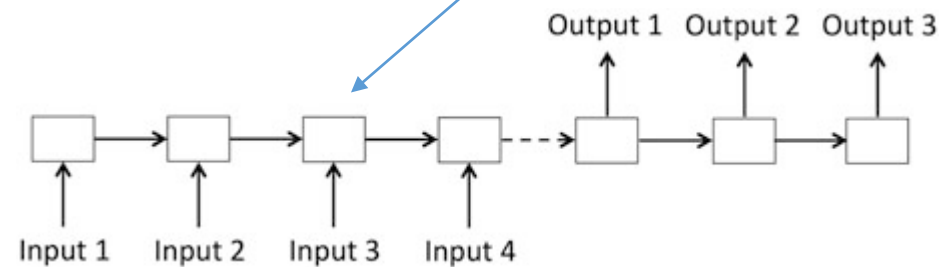
- Множество фаз тренировки.
- Различные типы ресурсов. Например, нужны фонемные словари.
- Большое количество гиперпараметров
- Различные функции оптимизации

End-To-End архитектуры для распознавания

- Connectionist Temporal Classification (CTC)



- Парадигма “Энкодер-декодер”



Рекуррентная формула

Сеть подает на вход в момент времени t выход $t-1$

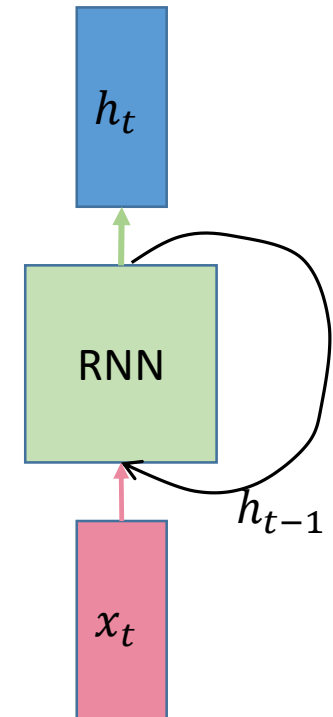
$$h_t = f_W(h_{t-1}, x_t)$$

Новое состояние

Старое состояние

Вход на шаге t

Функция с параметрами W



Рекуррентная формула

Сеть подает на вход в момент времени t выход $t-1$

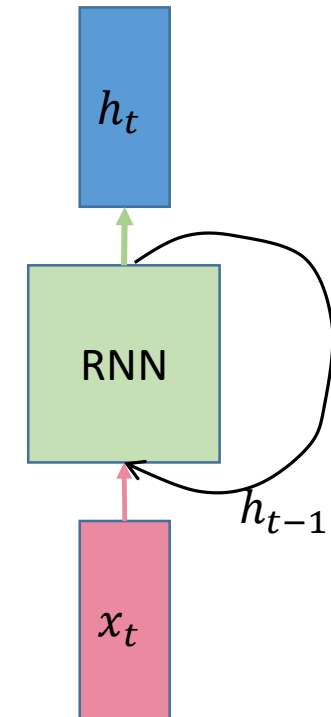
$$h_t = f_W(h_{t-1}, x_t)$$

Новое состояние

Старое состояние

Вход на шаге t

Функция с параметрами W



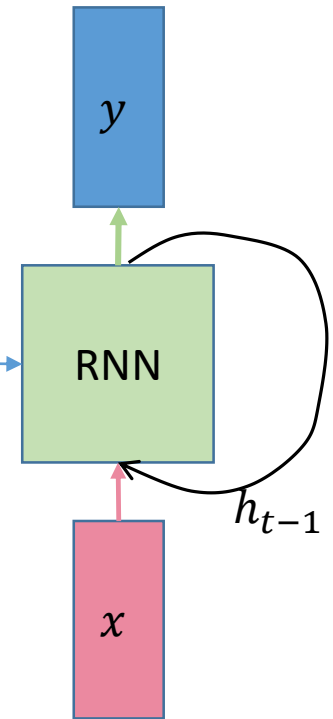
Для каждого момента времени используется одна функция и одна матрица весов

Базовая рекуррентная сеть

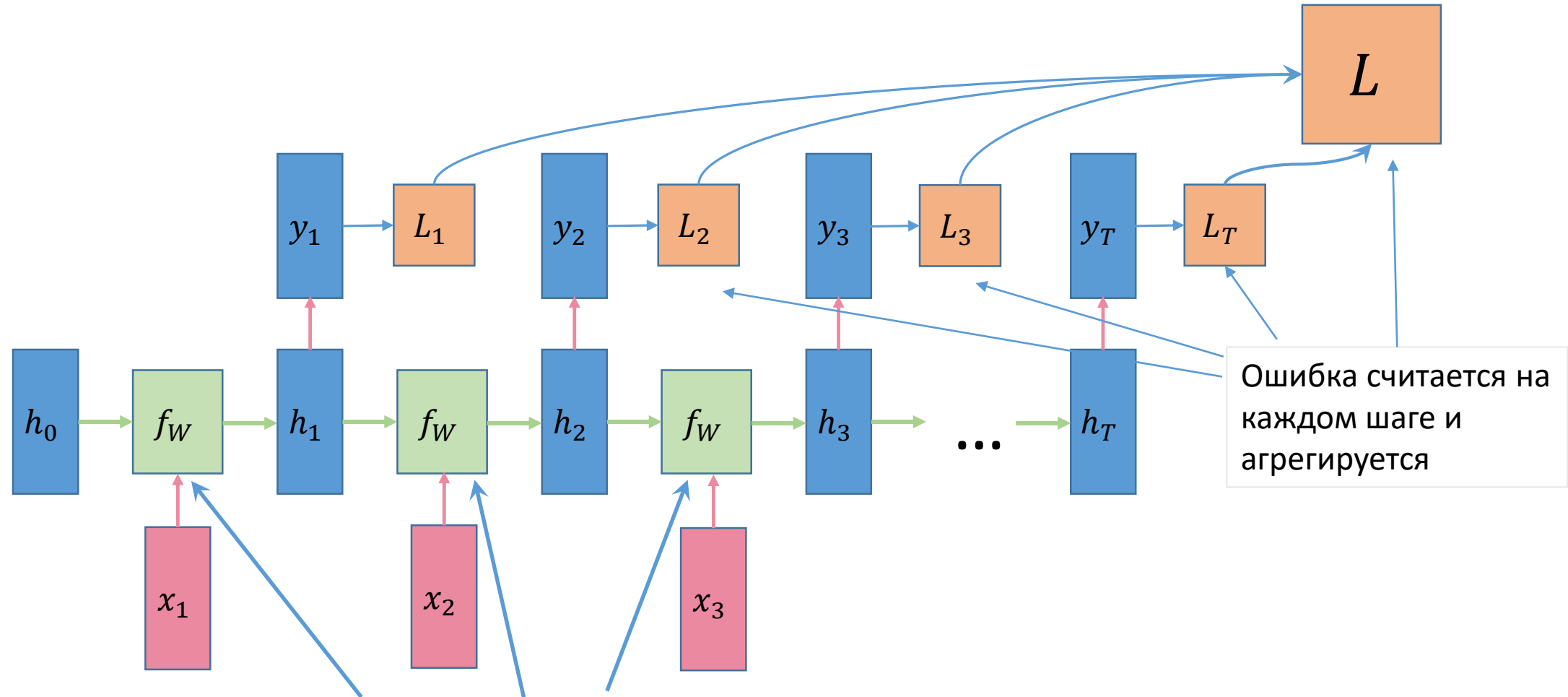
$$y_t = W_{hy} h_t$$

$$h_t = \tanh(W_{hh} h_{t-1} + W_{xh} x_t)$$

$$h_t = f_W(h_{t-1}, x_t)$$

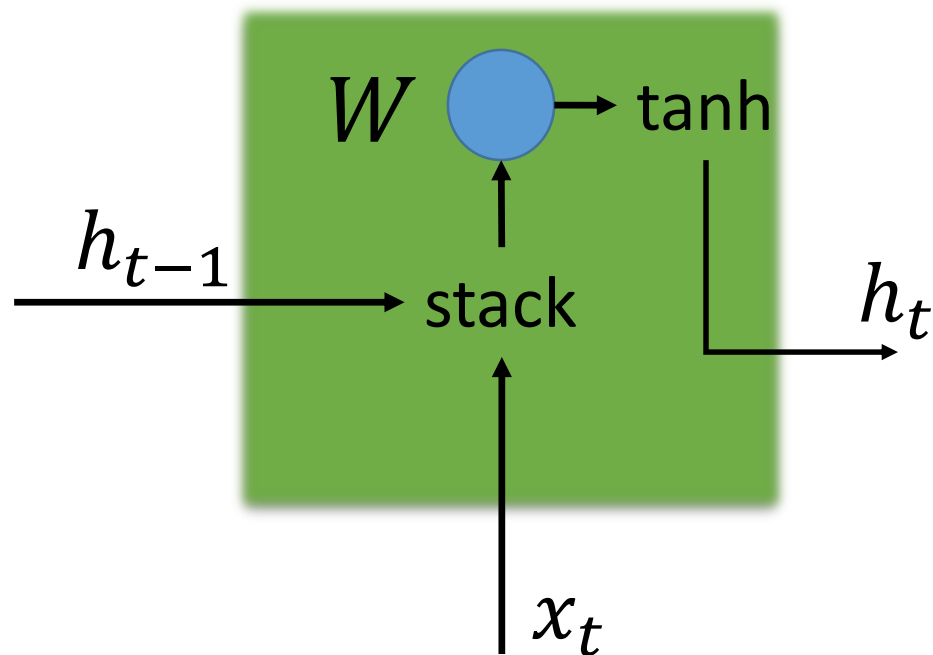


Вычислительный граф. Один ко многим



Одна матрица W на каждом шаге

Градиент через RNN

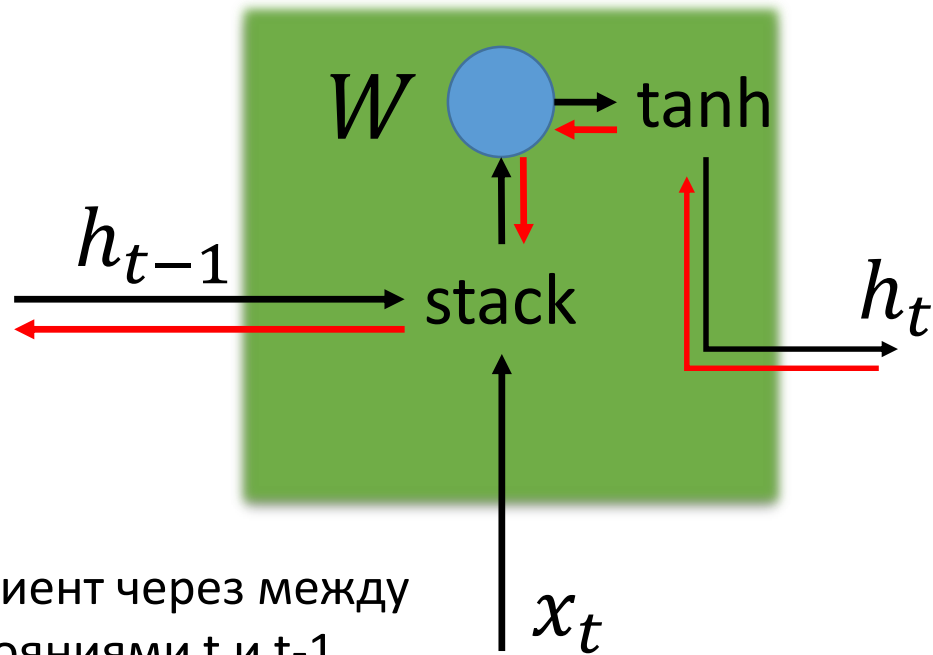


$$\begin{aligned} h_t &= \tanh(W_{hh}h_{t-1} + W_{xh}x_t) \\ &= \tanh\left(\begin{pmatrix} W_{hh} & W_{xh} \end{pmatrix} \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right) \\ &= \tanh\left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right) \end{aligned}$$

Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994

Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013

Градиент через RNN

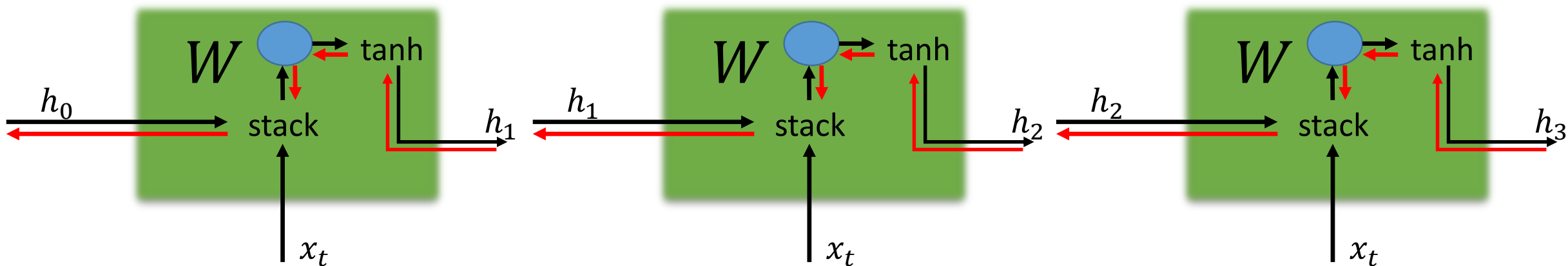


Градиент через между состояниями t и $t-1$ умножается на матрицу весов W^T

$$\begin{aligned} h_t &= \tanh(W_{hh}h_{t-1} + W_{xh}x_t) \\ &= \tanh\left(\begin{pmatrix} W_{hh} & W_{xh} \end{pmatrix} \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right) \\ &= \tanh\left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right) \end{aligned}$$

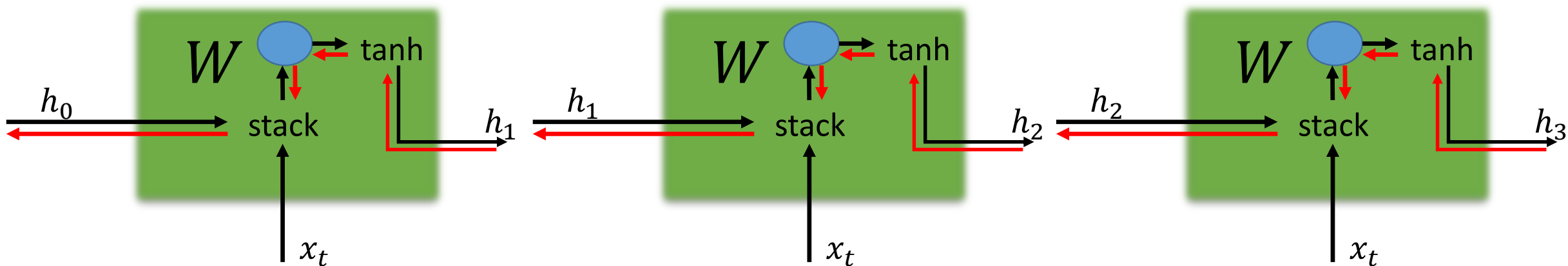
Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013

Градиент через RNN



Вычисление градиента для состояния 0 будет приводить к многократному умножению на матрицу весов W

Градиент через RNN



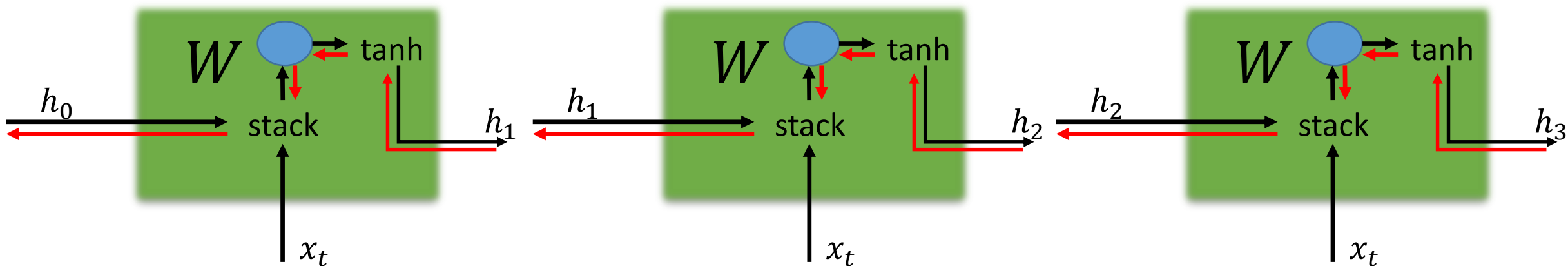
Вычисление градиента для состояния 0 будет приводить к многократному умножению на матрицу весов W

Значения матрицы больше 1:
Exploding gradient



Gradient clipping

Градиент через RNN



Вычисление градиента для состояния 0 будет приводить к многократному умножению на матрицу весов W

Значения матрицы меньше 1:
Vanishing gradient



Нужно менять архитектуру RNN

Long Short Term Memory (LSTM)

Vanilla RNN

$$h_t = \tanh \left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right)$$

LSTM

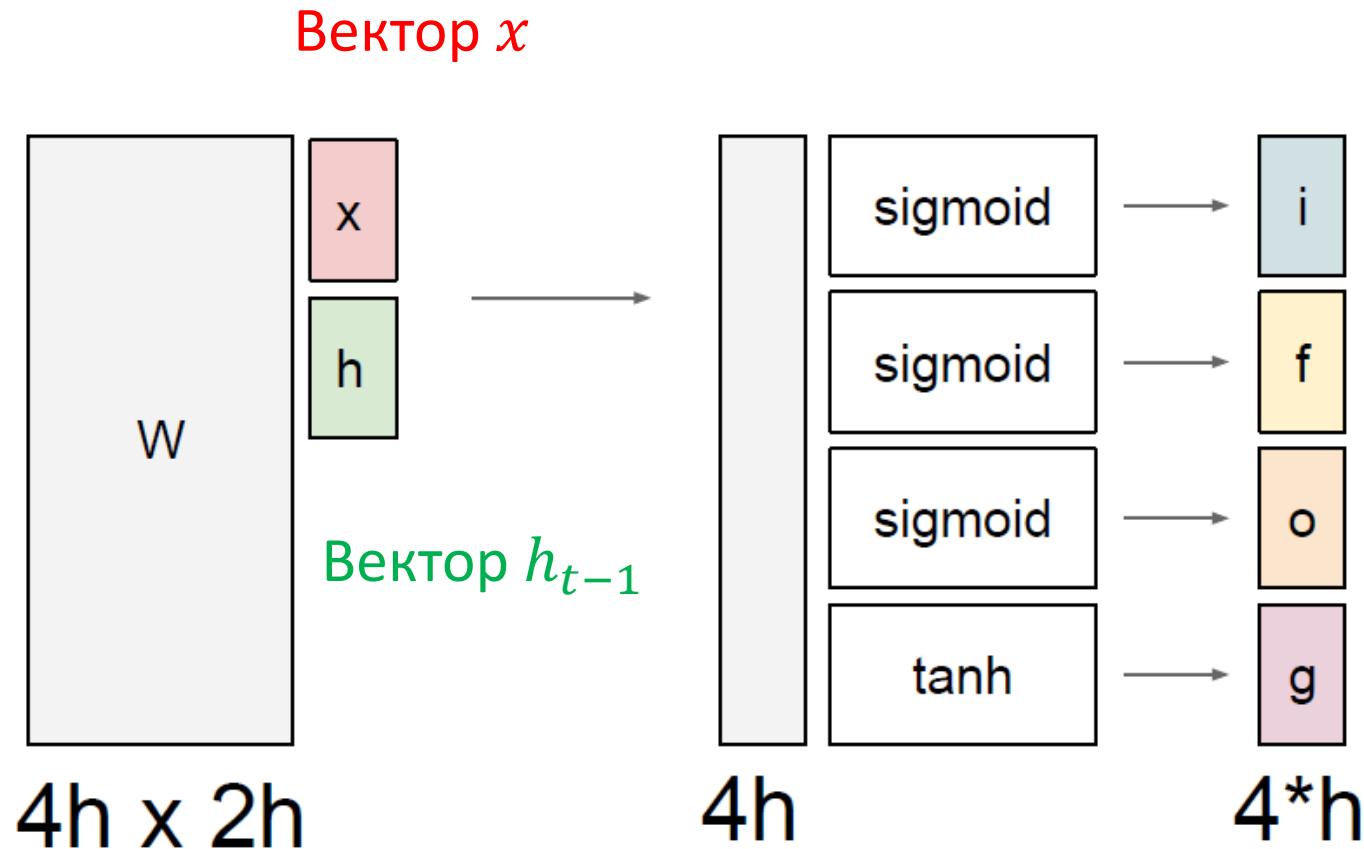
$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

Long Short Term Memory (LSTM)

f: Forget gate, Забывать ли состояние элемента
i: Input gate, Писать ли состояние элемента
g: Gate gate, Сколько писать в элемент
o: Output gate, Сколько вывести из элемента

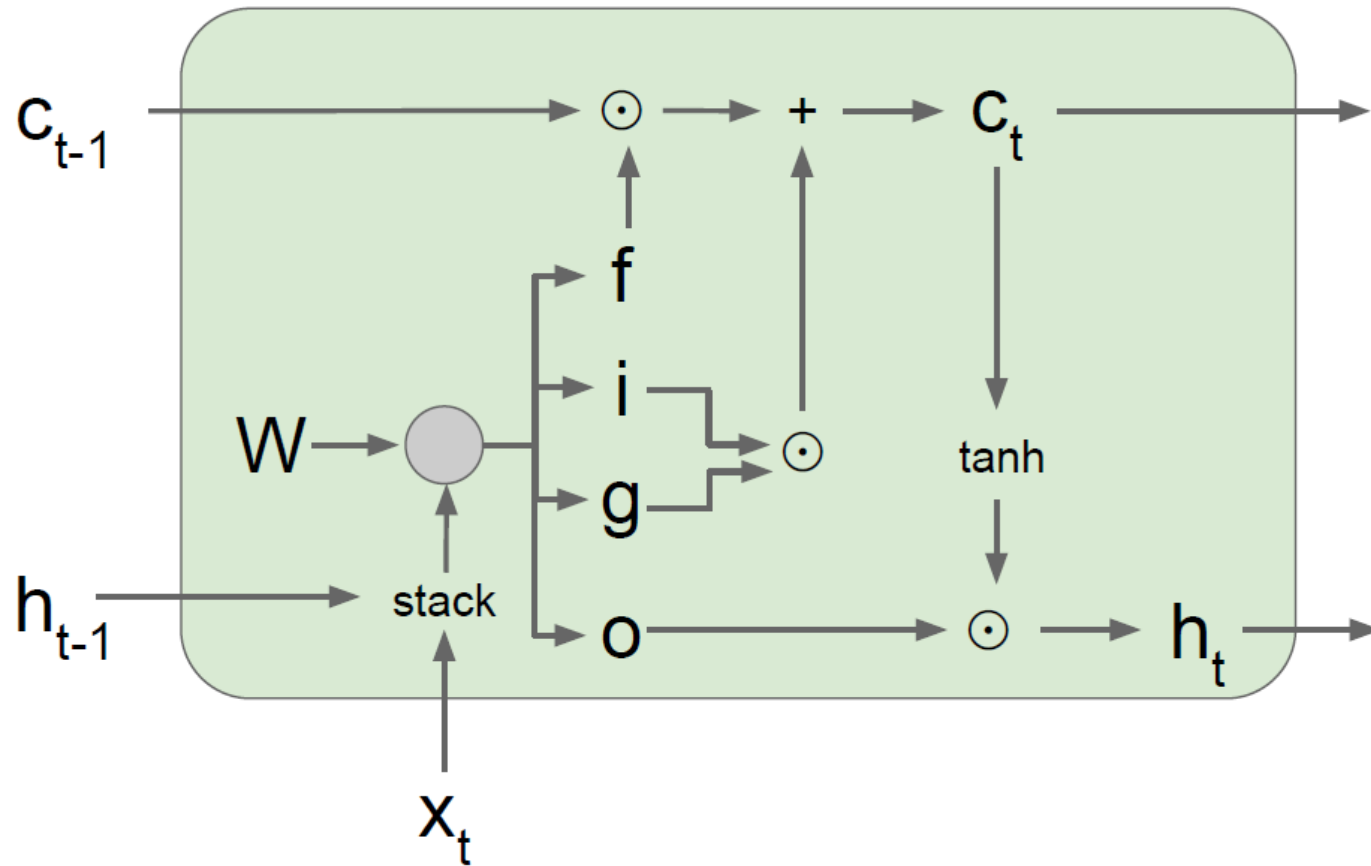


$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

Long Short Term Memory (LSTM)

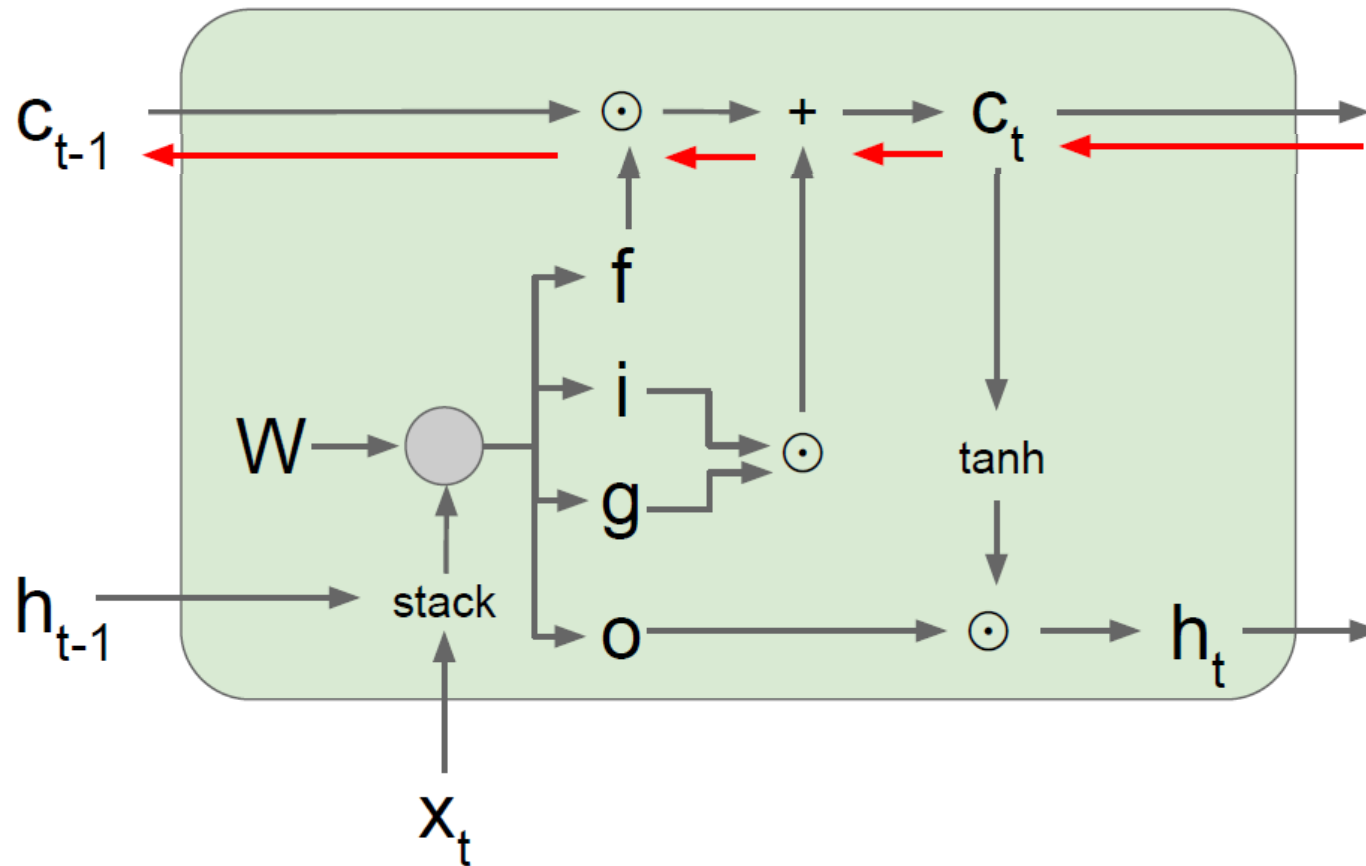


$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

LSTM. Gradient Flow



Обратный проход от c_t
до c_{t-1} перемножаются
только с f без матрицы W

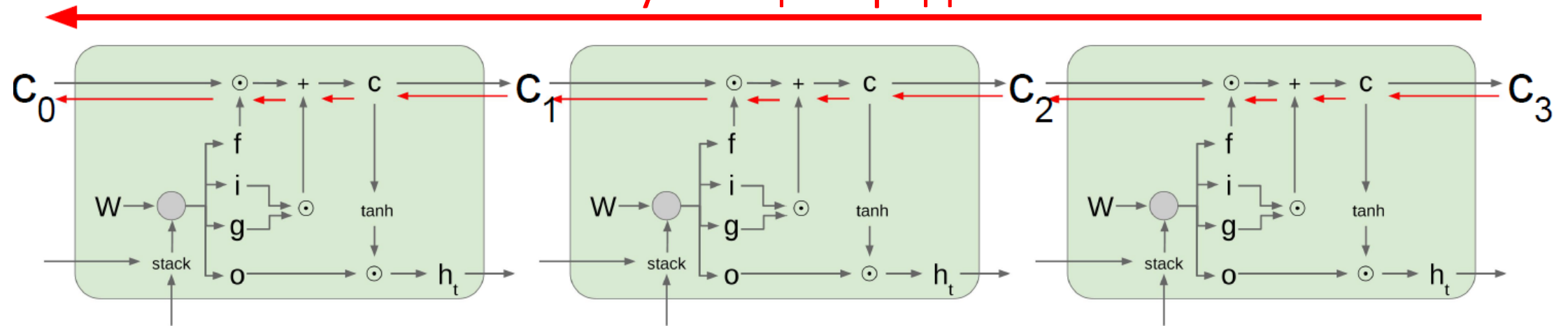
$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

LSTM. Gradient Flow

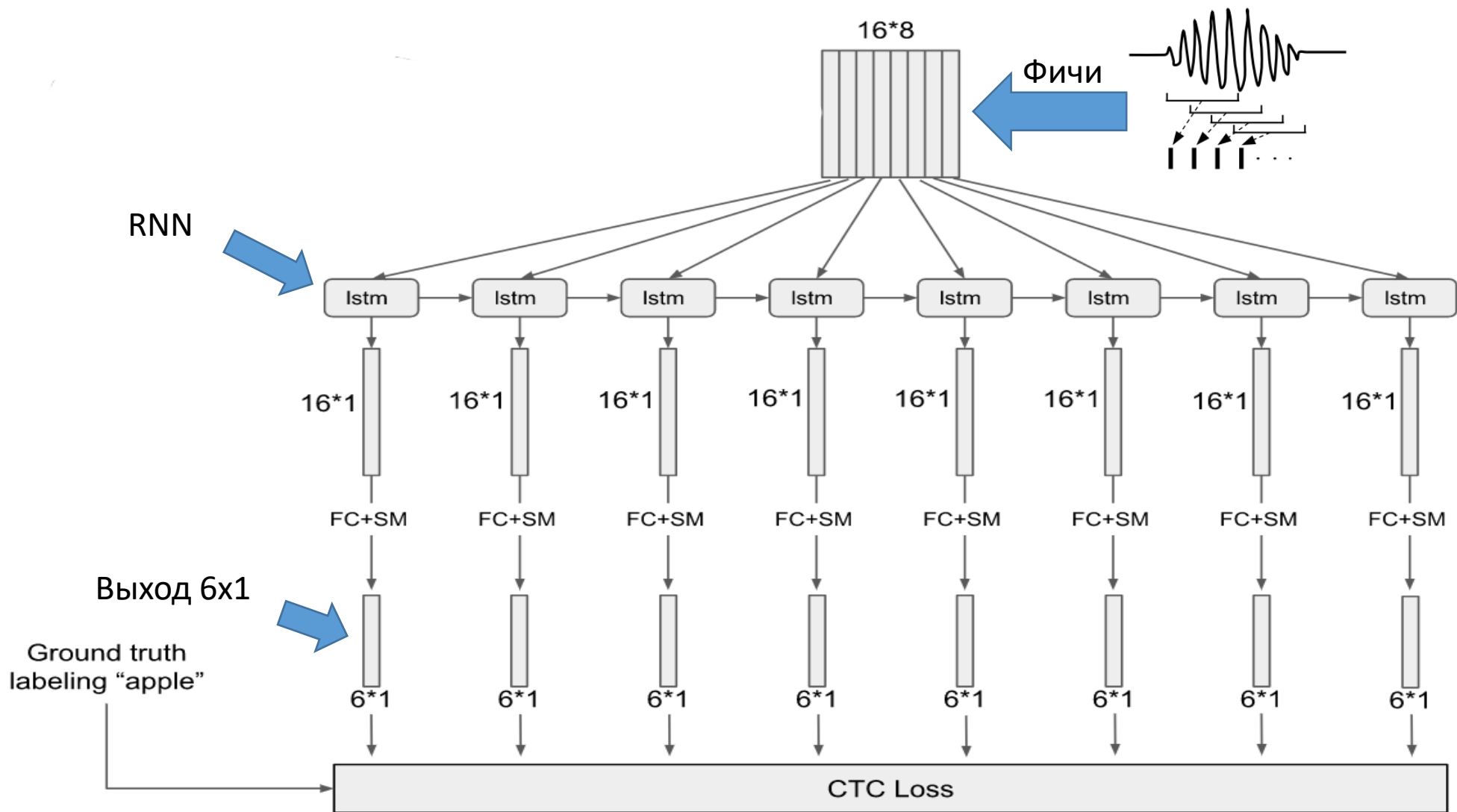
Незатухающий градиент



Connectionist Temporal Classification (CTC)

- Проблема RNN
 - Длина входной последовательности больше либо равна длине выходной последовательности
 - $\mathbf{x} = (x_1, x_2, \dots, x_T)$ – входная последовательность длиной T
 - $\mathbf{z} = (z_1, z_2, \dots, z_U)$ – выходная последовательность длиной U
 - $U \leq T$
 - Целевые метки берутся из алфавита L
 - Нет априорного способа выровнять \mathbf{x} и \mathbf{z}
- Наша цель использовать тренировочное множество S , чтобы построить $h : X \rightarrow Z$

CTC loss



CTC loss

- CTC loss – это “softmax” слой $p_l = \frac{\exp(x_l)}{\sum_k \exp(x_k)}$
- Количество выходов слоя на 1 больше, чем всего маркеров L
- Активация первых $|L|$ элементов слоя интерпретируется как вероятность
- Активация дополнительного юнита интерпретируется как отсутствие маркера. “blank”

CTC loss

- Для входной последовательности x длиной T
 - Задаем RNN с m входами, n выходами и w – вектор весов как непрерывное отображение $N_w: (R_m)^T \rightarrow (R_n)^T$
 - Тогда $y = N_w(x)$ – последовательность выходов RNN
- y_k^t - активация выходного элемента k в момент времени t
- y_k^t - вероятность пронаблюдать маркер k в момент времени t
 - Определяет распределение по множеству L'^T последовательностей длины T над алфавитом $L' = L \cup \{blank\}$:
 - $p(\pi|x) = \prod_{t=1}^T y_{\pi_t}^t, \forall \pi \in L'^T$.
- Элементы L'^T - это пути π

CTC loss - пути



Path1: "ap-pl-ee" $\xrightarrow{B(\text{"ap-pl-ee"})}$ Labeling: "apple"

$$p(\text{"ap-pl-ee"}) = y_a^1 \cdot y_p^2 \cdot y_{-}^3 \cdot y_p^4 \cdot y_l^5 \cdot y_{-}^6 \cdot y_e^7 \cdot y_e^8$$

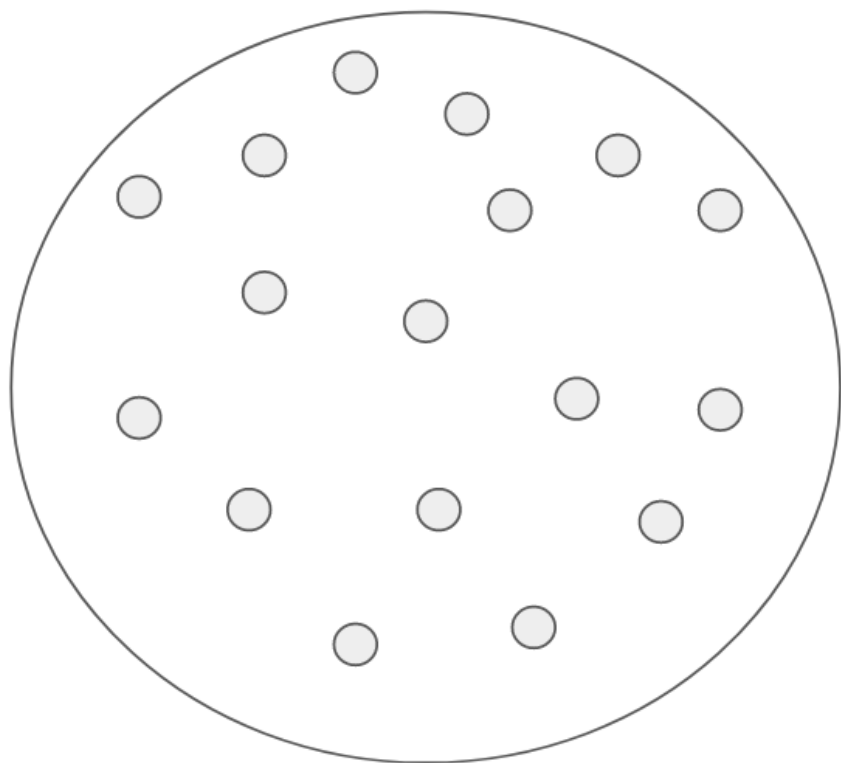
CTC loss - пути



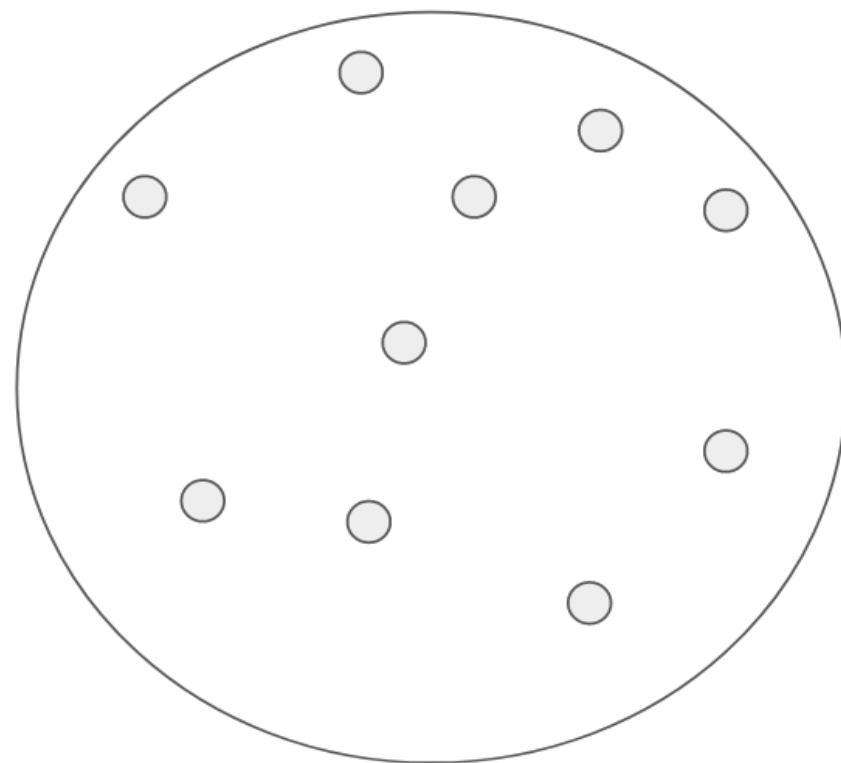
Path1: "ap-pl-ee" $\xrightarrow{B("ap-pl-ee")}$ Labeling: "apple"
 $p("ap-pl-ee") = y_a^1 \cdot y_p^2 \cdot y_-^3 \cdot y_p^4 \cdot y_l^5 \cdot y_-^6 \cdot y_e^7 \cdot y_e^8$

Path2: "aapp--le" $\xrightarrow{B("aapp--le")}$ Labeling: "apple"
 $p("aapp--le") = y_a^1 \cdot y_a^2 \cdot y_p^3 \cdot y_p^4 \cdot y_-^5 \cdot y_-^6 \cdot y_l^7 \cdot y_e^8$

CTC loss. Пути



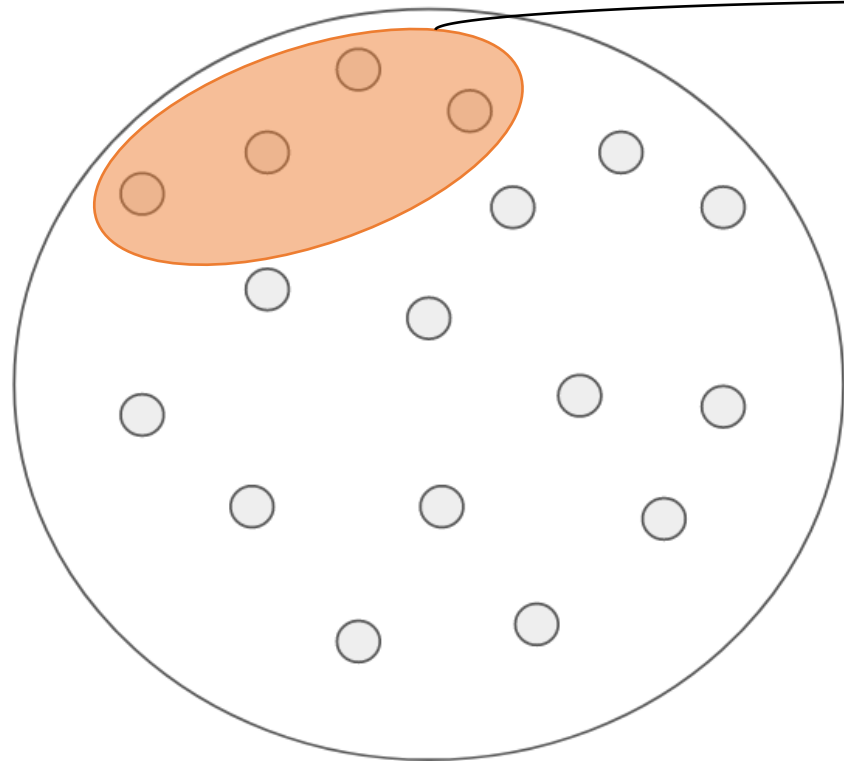
Set of all possible paths



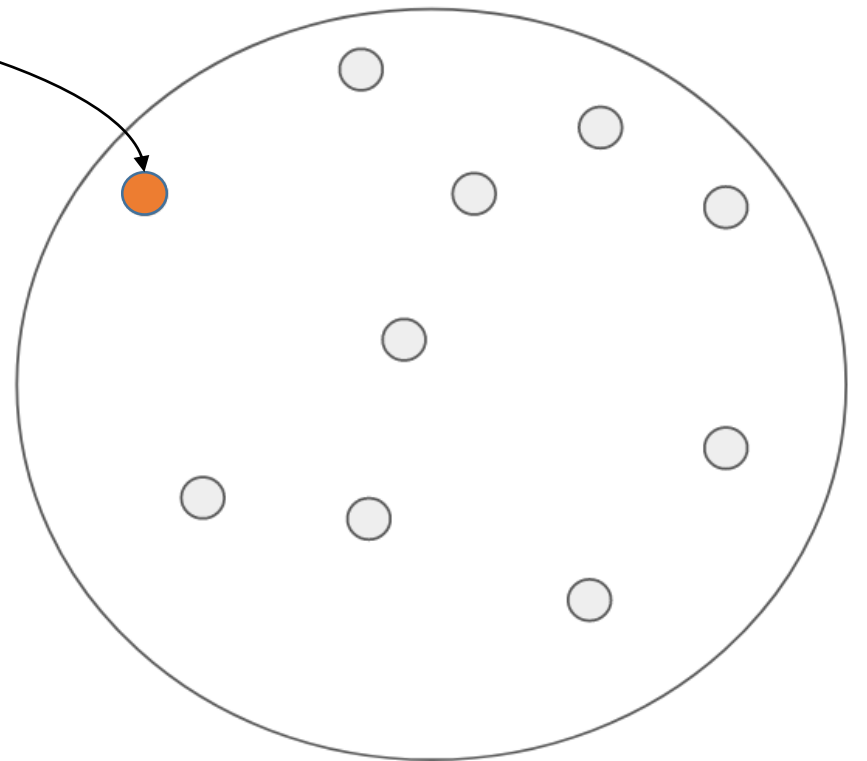
Set of all possible labelings

CTC loss. Пути

Пути, соответствующие метке



Set of all possible paths



Set of all possible labelings

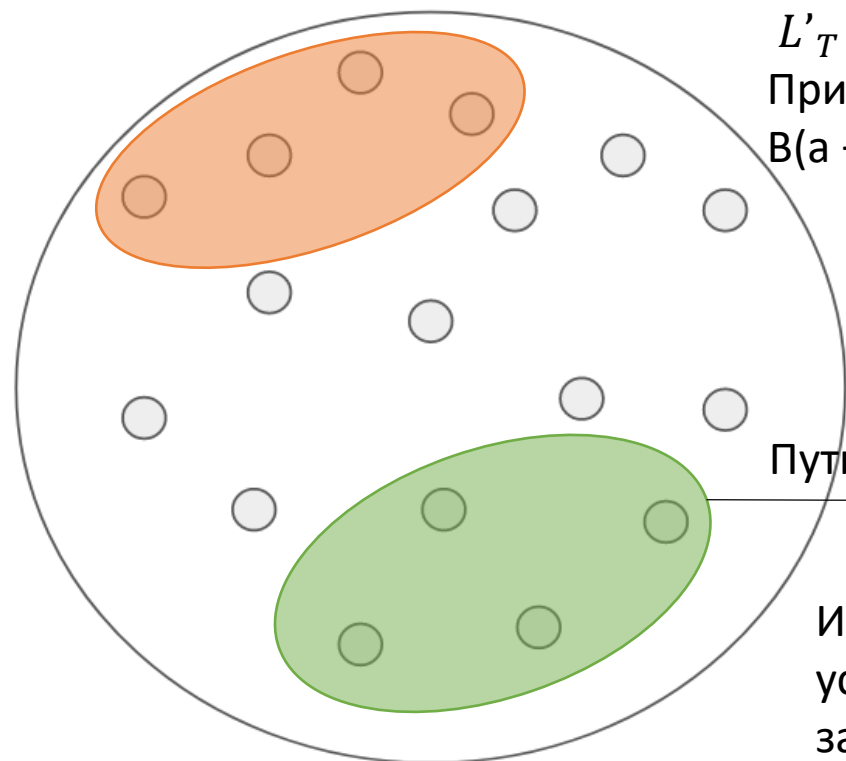
CTC loss. Пути

Мы хотим получить маппинг: $B :$

$$L^T \rightarrow L^{\leq T}.$$

Пример:

$$B(a - ab-) = B(-aa - -abb) = aab).$$

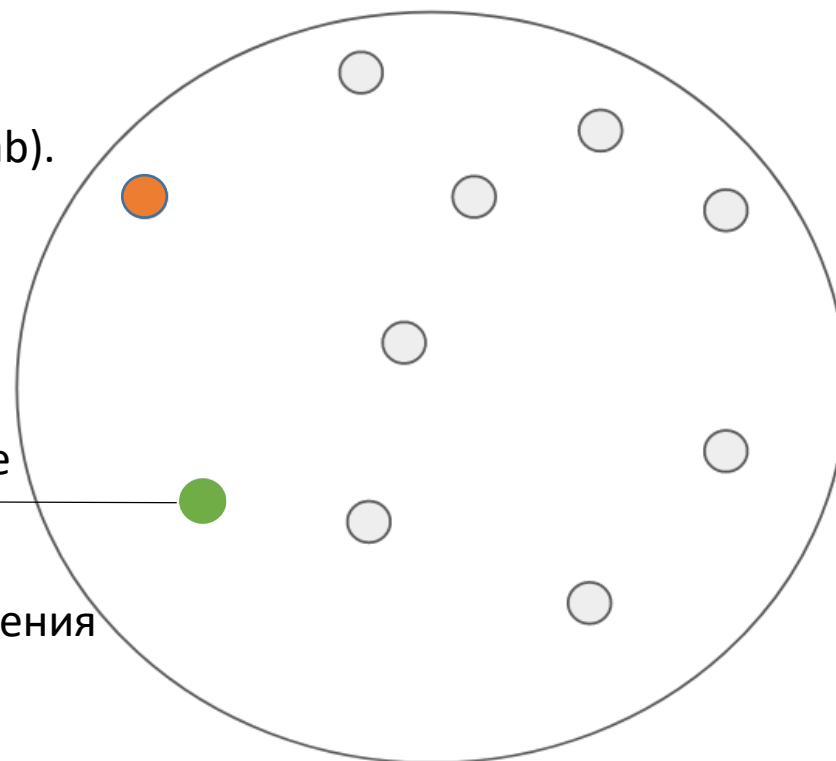


Set of all possible paths

Пути, соответствующие метке

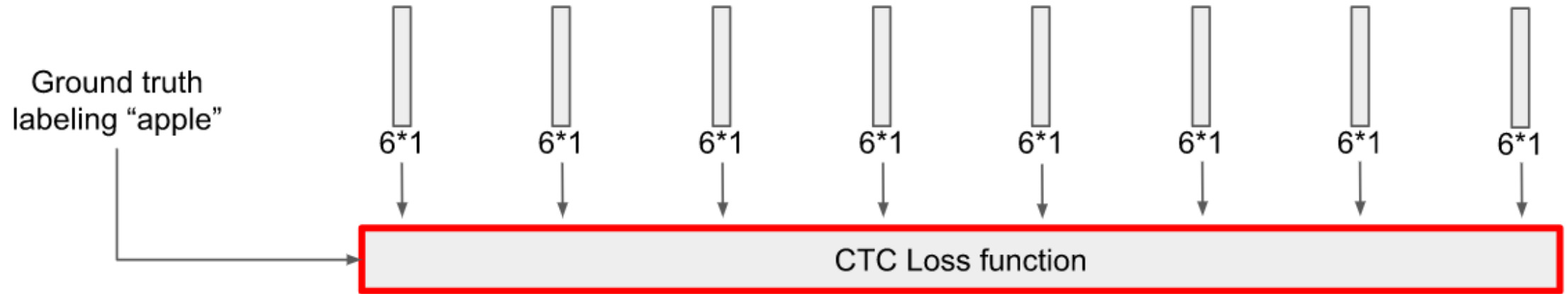
Используем B для определения условной вероятности B заданной метки $l \in L^{\leq T}$

$$p(l|x) = \sum_{X_{\pi} \in B^{-1}(l)} p(\pi|x).$$



Set of all possible labelings

CTC loss



$$\text{CTC Loss} = -\ln(p(\text{"apple"}))$$

- Вероятность слова – сумма вероятностей по всем возможным путям
- $6^8 = 1\,679\,616$ – возможных путей (случай из примера)
- Используем динамическое программирование для нахождения вероятности целевой последовательности

Нахождение возможных путей

- Аналогично прямому проходу и проходам в НММ мы рассчитываем α и β

$$\alpha_t(s) \stackrel{\text{def}}{=} \sum_{\substack{\pi \in N^T: \\ \mathcal{B}(\pi_{1:t}) = \mathbf{l}_{1:s}}} \prod_{t'=1}^t y_{\pi_{t'}}.$$

- Суммарная вероятность всех путей, чей префикс заканчивается символом в позиции s в момент времени t от начала последовательности

$$\beta_t(s) \stackrel{\text{def}}{=} \sum_{\substack{\pi \in N^T: \\ \mathcal{B}(\pi_{t:T}) = \mathbf{l}_{s:|l|}}} \prod_{t'=t}^T y_{\pi_{t'}}.$$

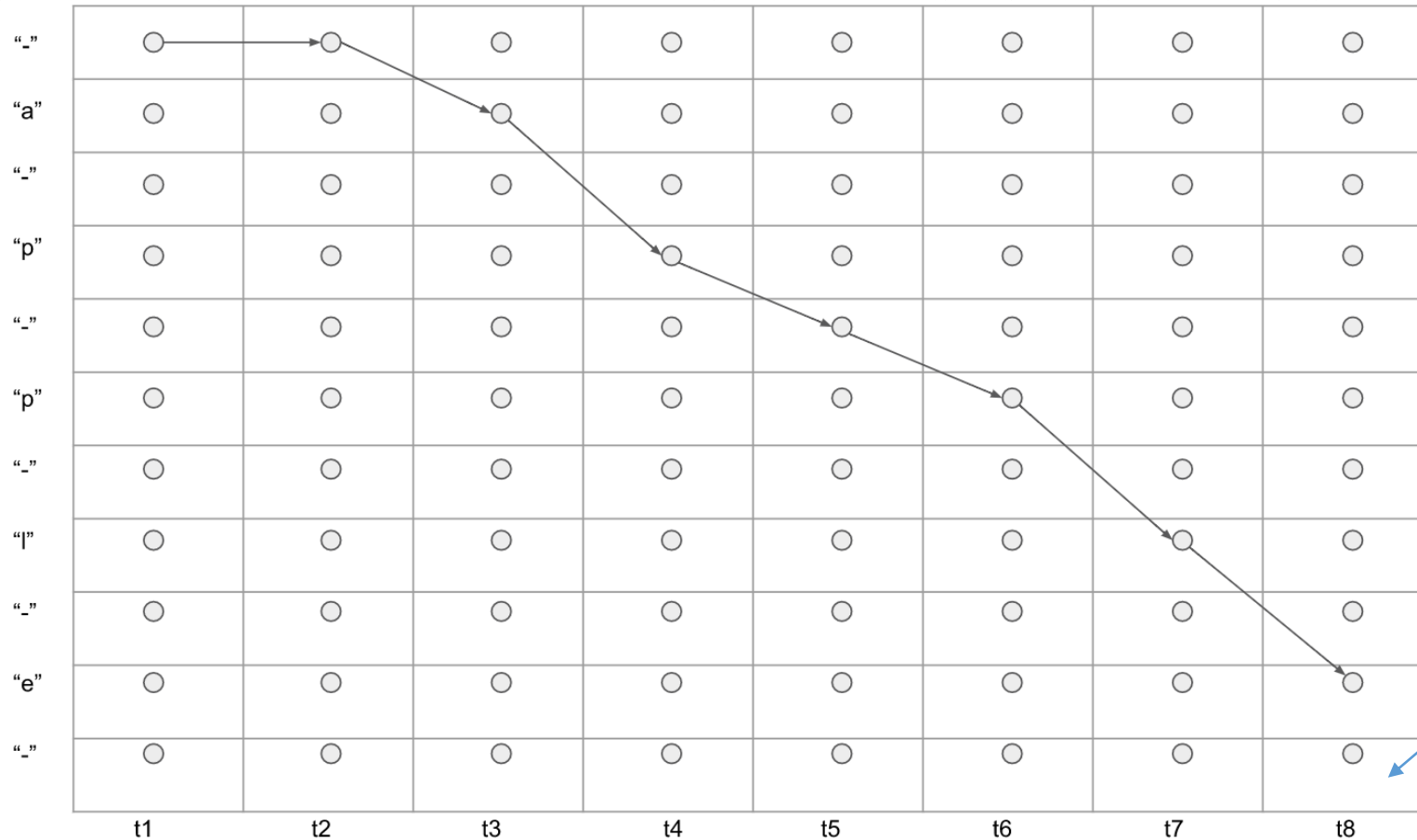
- Суммарная вероятность всех путей, чей суффикс начинается с символом в позиции s в момент времени t

Нахождение возможных путей

"_"	○	○	○	○	○	○	○	○
"a"	○	○	○	○	○	○	○	○
"_"	○	○	○	○	○	○	○	○
"p"	○	○	○	○	○	○	○	○
"_"	○	○	○	○	○	○	○	○
"p"	○	○	○	○	○	○	○	○
"_"	○	○	○	○	○	○	○	○
"l"	○	○	○	○	○	○	○	○
"_"	○	○	○	○	○	○	○	○
"e"	○	○	○	○	○	○	○	○
"_"	○	○	○	○	○	○	○	○
	t1	t2	t3	t4	t5	t6	t7	t8

Нахождение возможных путей

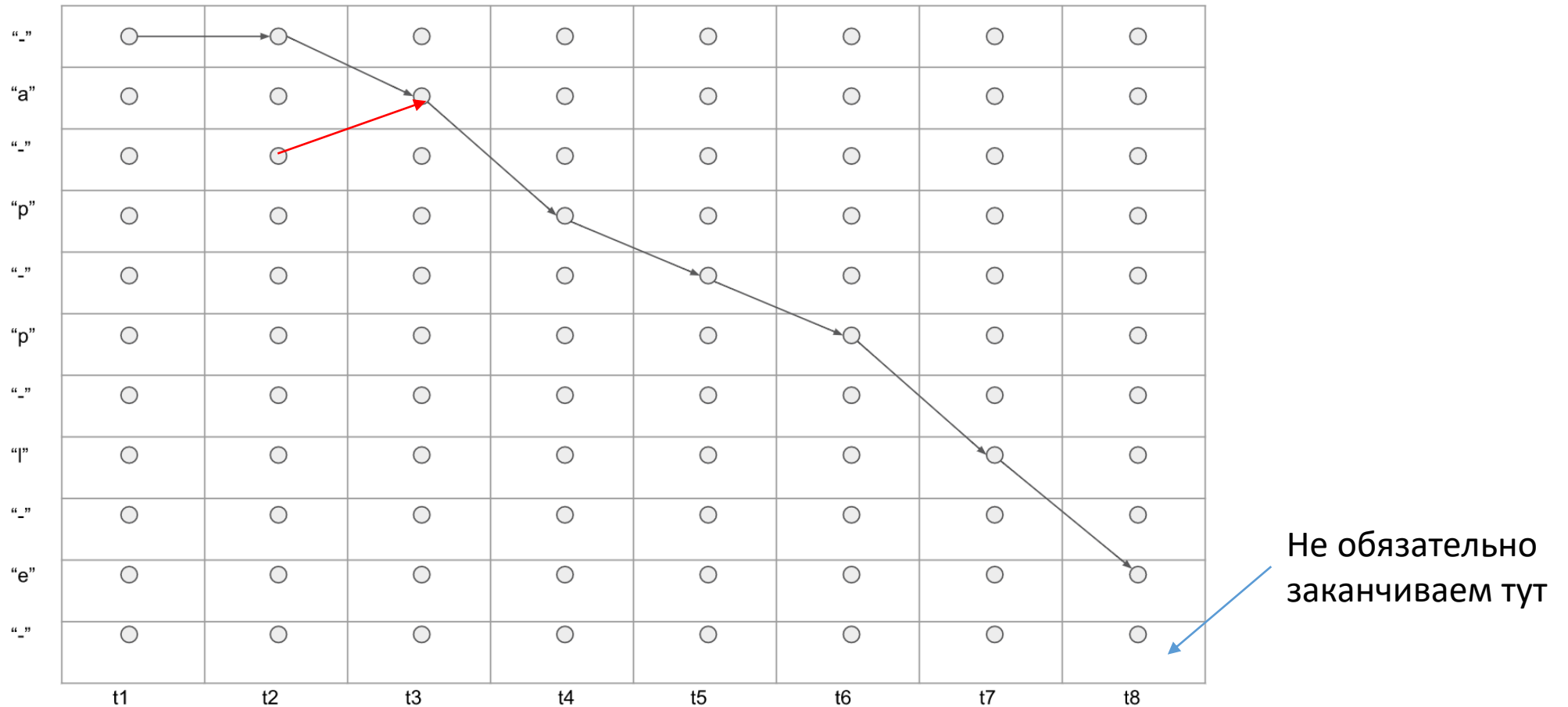
Пример, путь “_arple” может быть отмапирован на маркер “apple” $B(--ar-ple) = “apple”$



Не обязательно заканчиваем тут

Нахождение возможных путей

Пример: невозможный переход, нельзя предсказывать предыдущий символ.



Нахождение возможных путей

Начинаем либо b или с первого символа

Инициализация

$$\alpha_1(1) = y_b^1$$

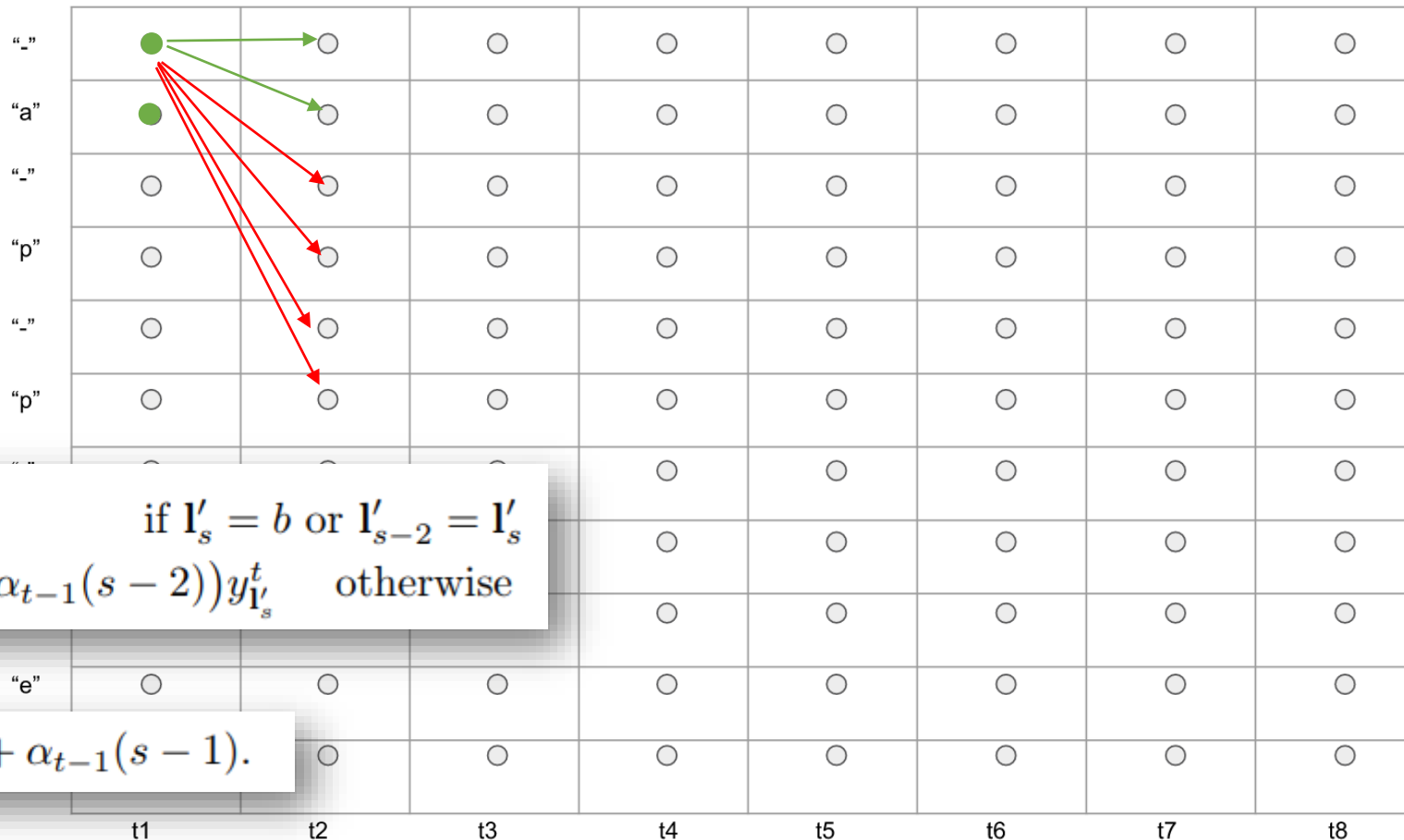
$$\alpha_1(2) = y_{I_1}^1$$

$$\alpha_1(s) = 0, \forall s > 2$$

"_"	○	○	○	○	○	○	○	○
"a"	○	○	○	○	○	○	○	○
"_"	○	○	○	○	○	○	○	○
"p"	○	○	○	○	○	○	○	○
"_"	○	○	○	○	○	○	○	○
"p"	○	○	○	○	○	○	○	○
"_"	○	○	○	○	○	○	○	○
"l"	○	○	○	○	○	○	○	○
"_"	○	○	○	○	○	○	○	○
"e"	○	○	○	○	○	○	○	○
"_"	○	○	○	○	○	○	○	○
	t1	t2	t3	t4	t5	t6	t7	t8

Нахождение возможных путей

Куда можно двигаться из этих точек

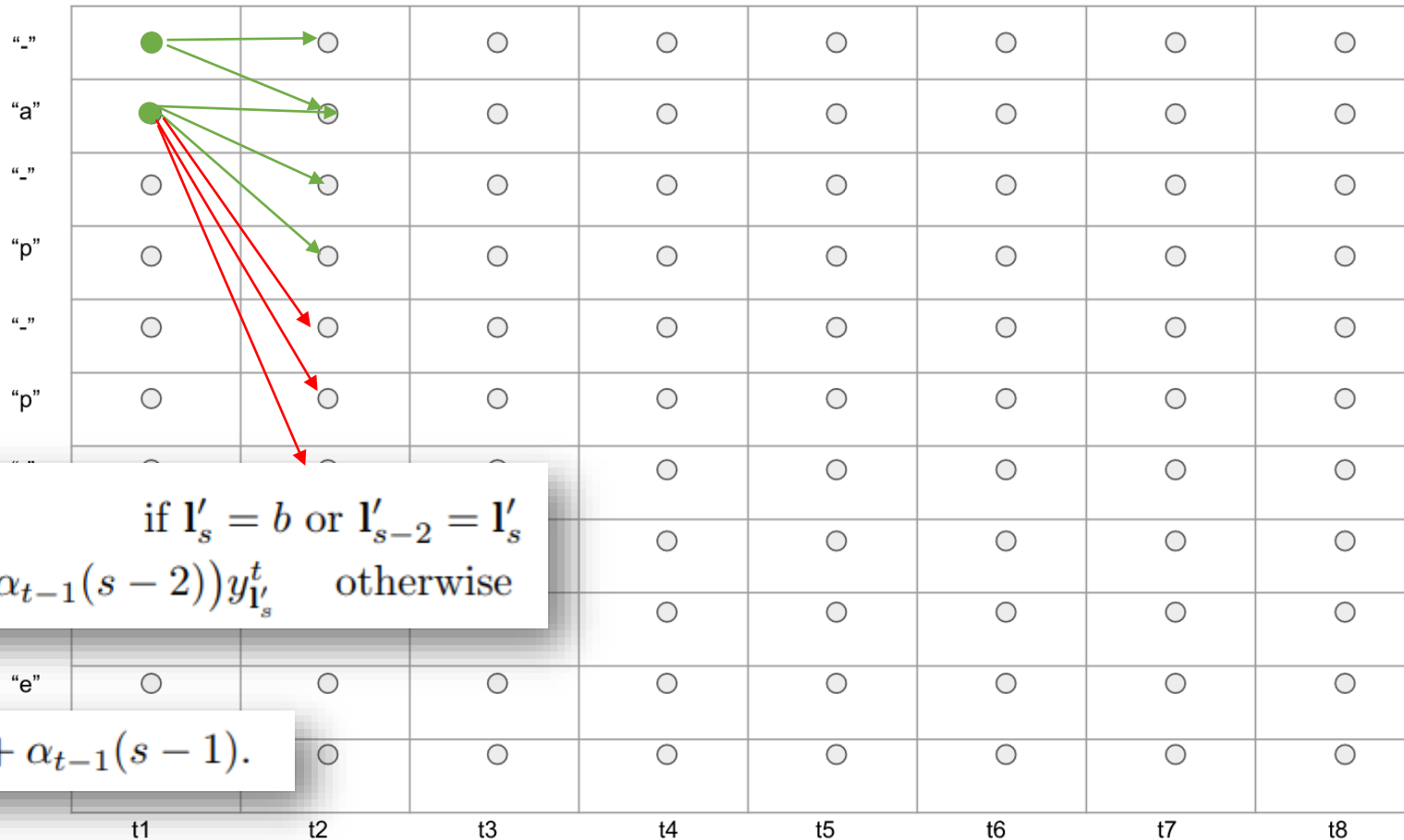


$$\alpha_t(s) = \begin{cases} \bar{\alpha}_t(s) y_{l'_s}^t & \text{if } l'_s = b \text{ or } l'_{s-2} = l'_s \\ (\bar{\alpha}_t(s) + \alpha_{t-1}(s-2)) y_{l'_s}^t & \text{otherwise} \end{cases}$$

$$\bar{\alpha}_t(s) \stackrel{\text{def}}{=} \alpha_{t-1}(s) + \alpha_{t-1}(s-1).$$

Нахождение возможных путей

Куда можно двигаться из этих точек

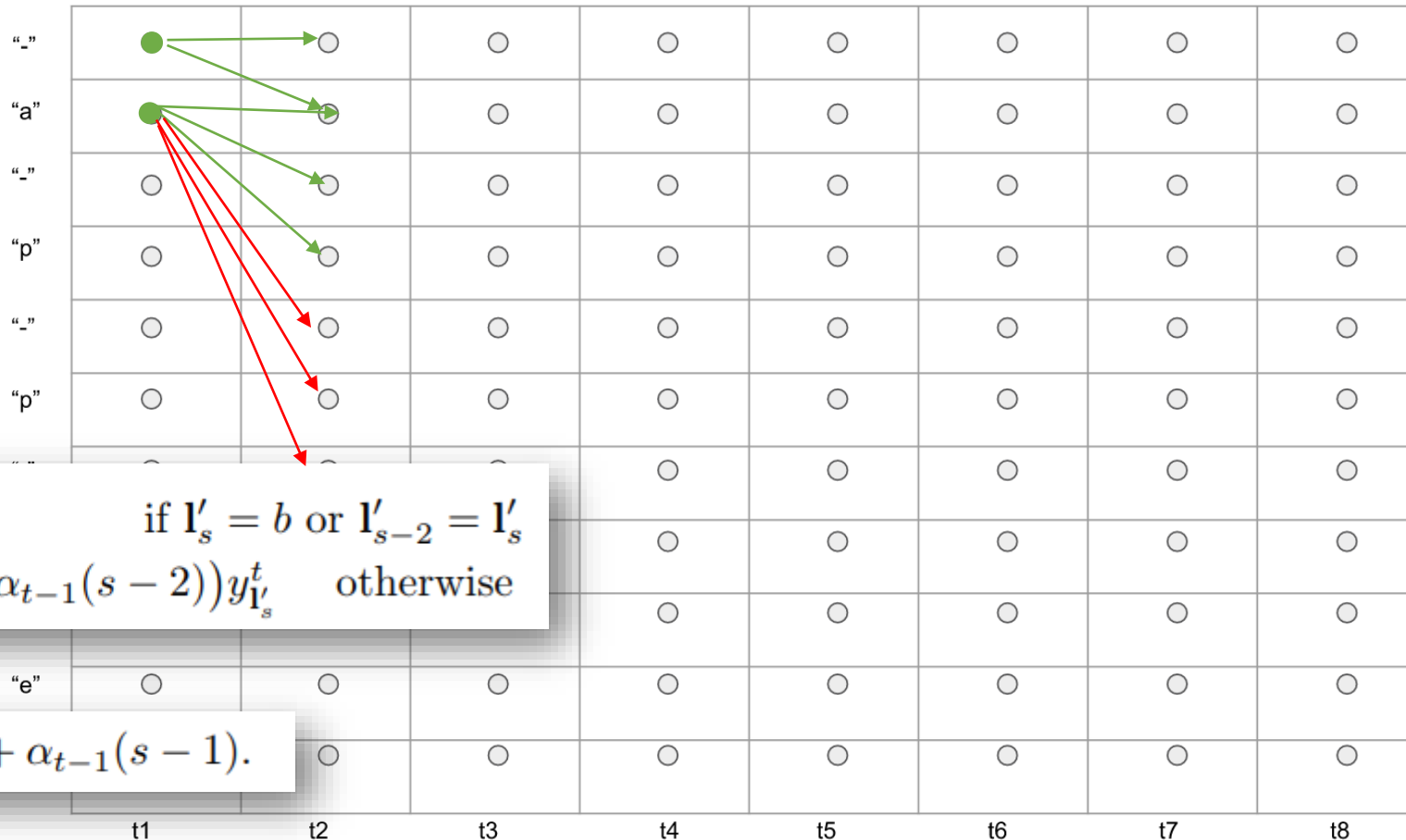


$$\alpha_t(s) = \begin{cases} \bar{\alpha}_t(s) y_{l'_s}^t & \text{if } l'_s = b \text{ or } l'_{s-2} = l'_s \\ (\bar{\alpha}_t(s) + \alpha_{t-1}(s-2)) y_{l'_s}^t & \text{otherwise} \end{cases}$$

$$\bar{\alpha}_t(s) \stackrel{\text{def}}{=} \alpha_{t-1}(s) + \alpha_{t-1}(s-1).$$

Нахождение возможных путей

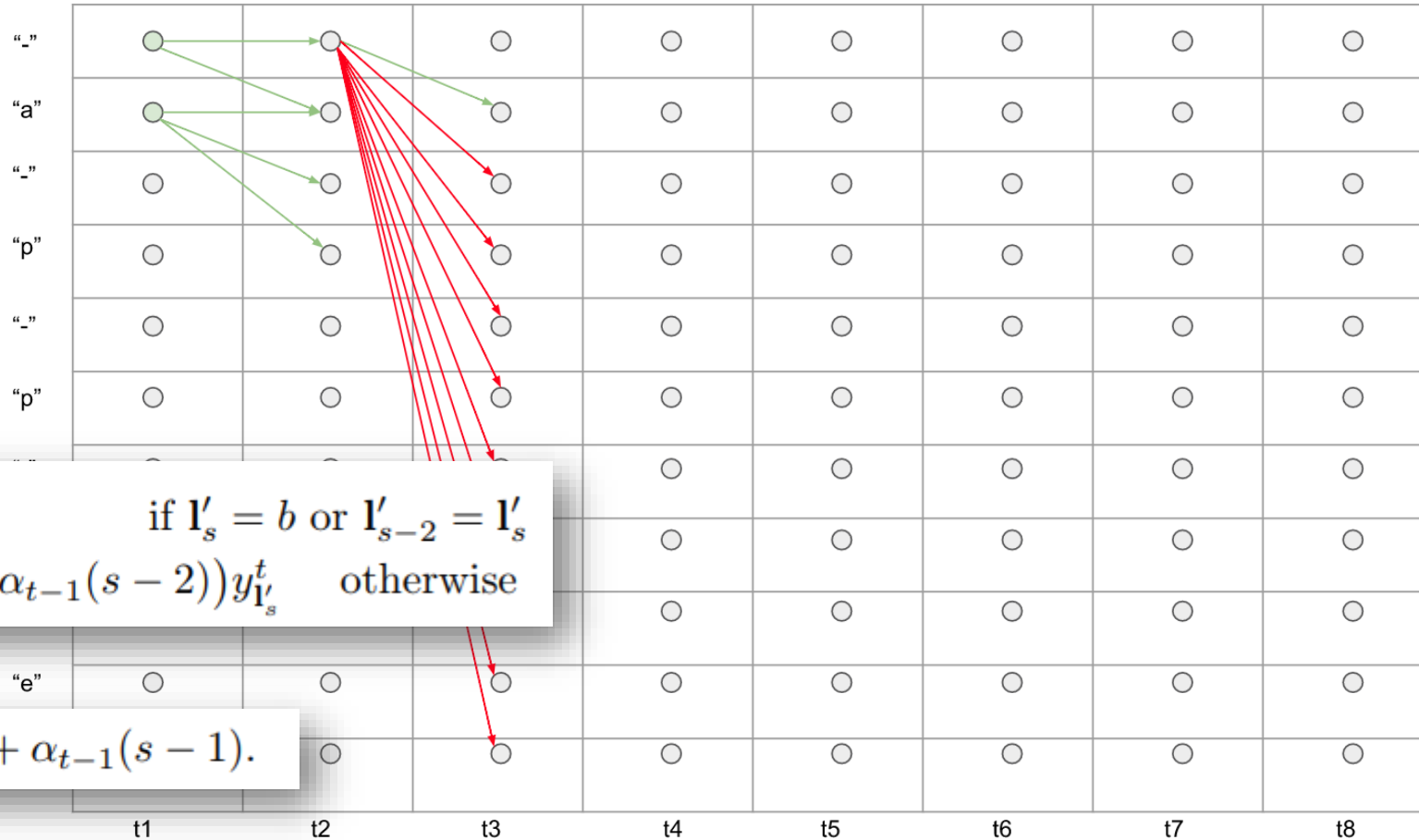
Куда можно двигаться из этих точек



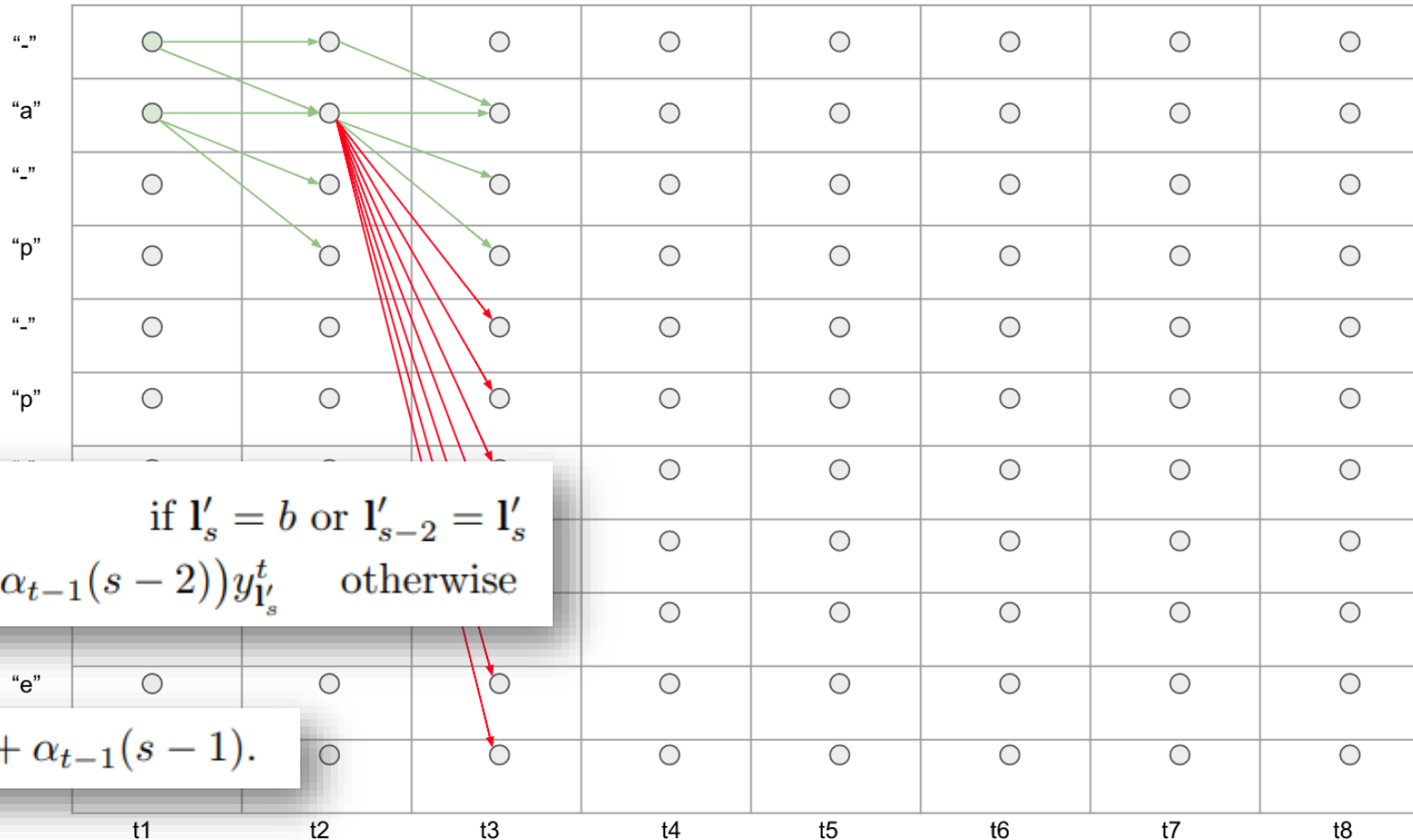
$$\alpha_t(s) = \begin{cases} \bar{\alpha}_t(s) y_{l'_s}^t & \text{if } l'_s = b \text{ or } l'_{s-2} = l'_s \\ (\bar{\alpha}_t(s) + \alpha_{t-1}(s-2)) y_{l'_s}^t & \text{otherwise} \end{cases}$$

$$\bar{\alpha}_t(s) \stackrel{\text{def}}{=} \alpha_{t-1}(s) + \alpha_{t-1}(s-1).$$

Нахождение возможных путей



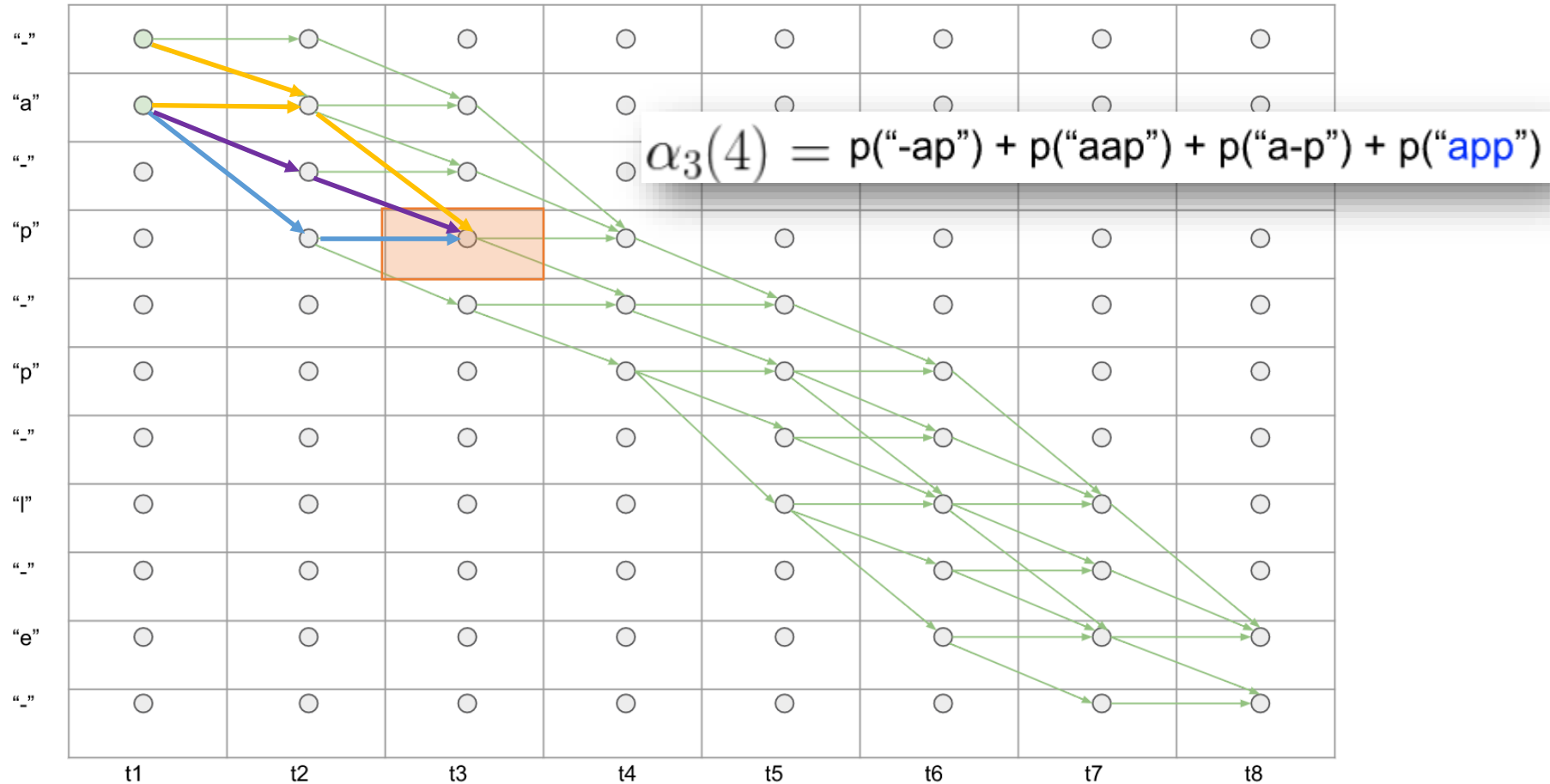
Нахождение возможных путей



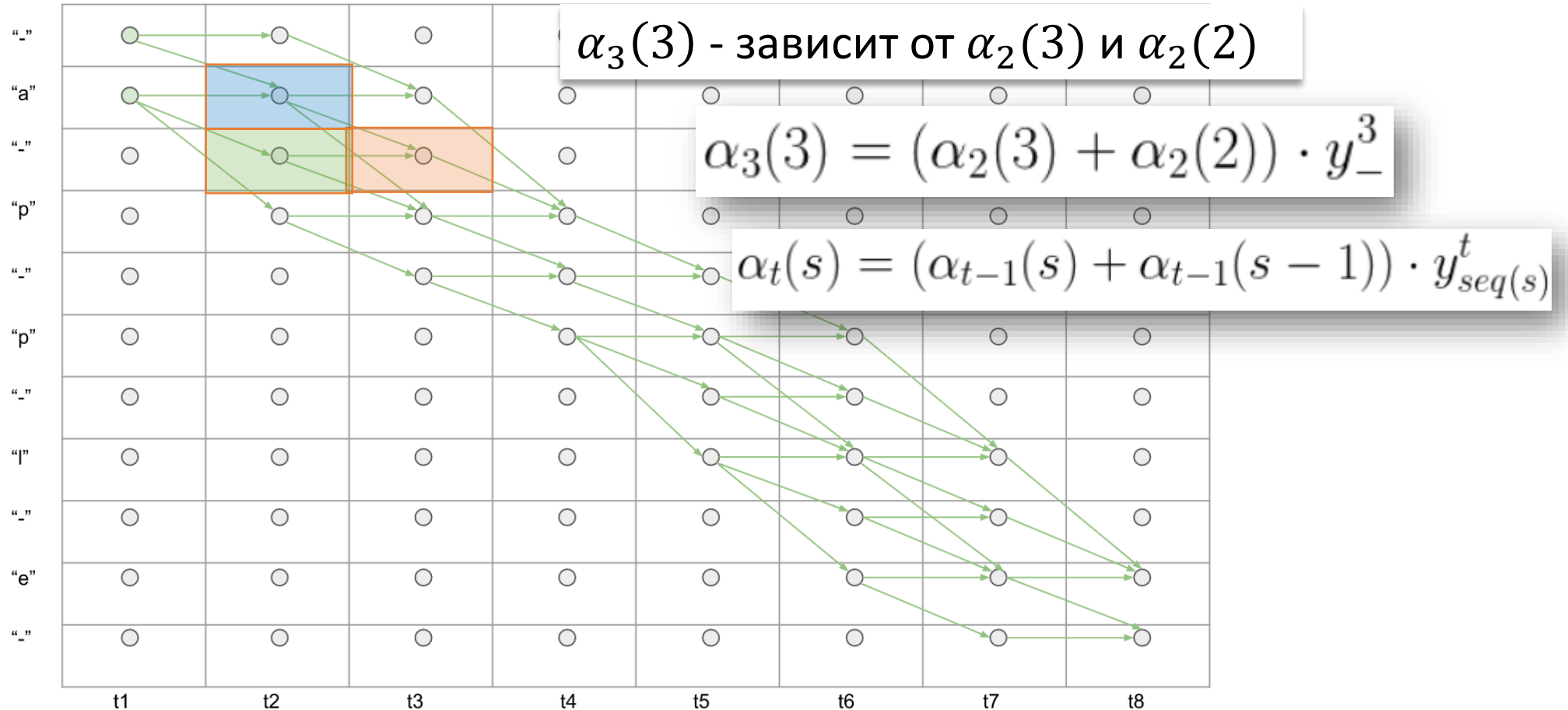
$$\alpha_t(s) = \begin{cases} \bar{\alpha}_t(s) y_{l'_s}^t & \text{if } l'_s = b \text{ or } l'_{s-2} = l'_s \\ (\bar{\alpha}_t(s) + \alpha_{t-1}(s-2)) y_{l'_s}^t & \text{otherwise} \end{cases}$$

$$\bar{\alpha}_t(s) \stackrel{\text{def}}{=} \alpha_{t-1}(s) + \alpha_{t-1}(s-1).$$

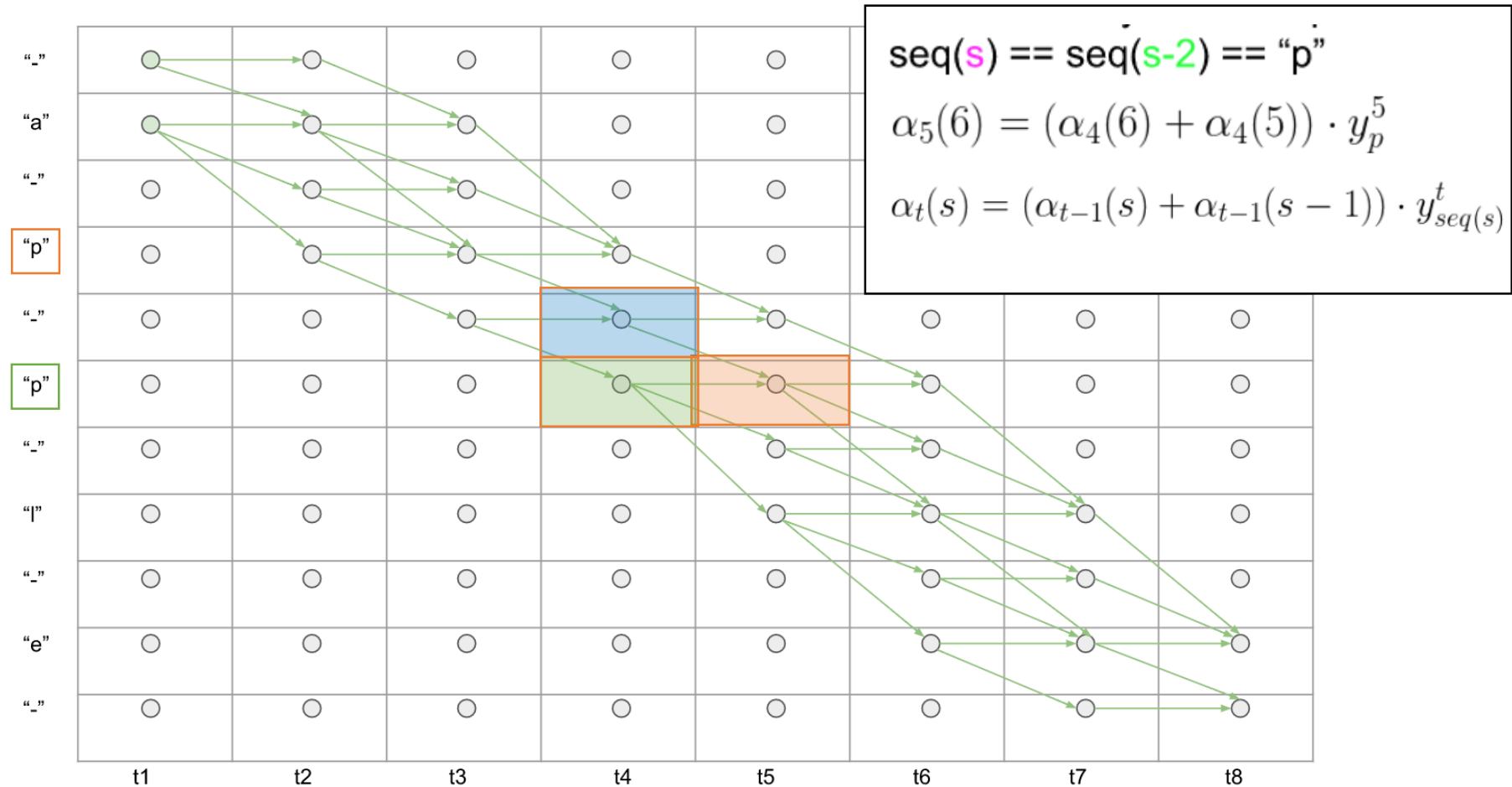
Нахождение возможных путей



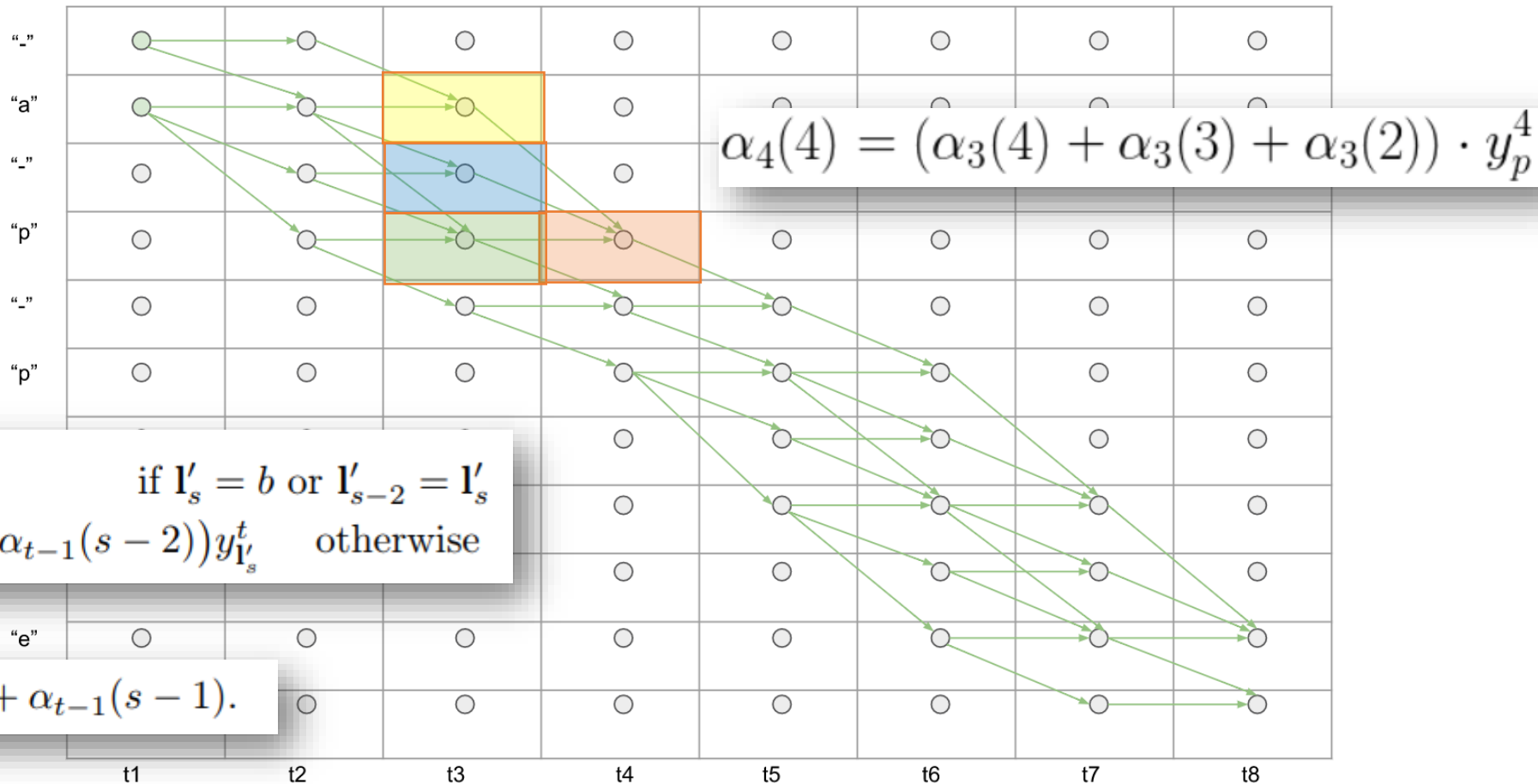
Нахождение возможных путей



Нахождение возможных путей



Нахождение возможных путей



Нахождение возможных путей. Обратный проход

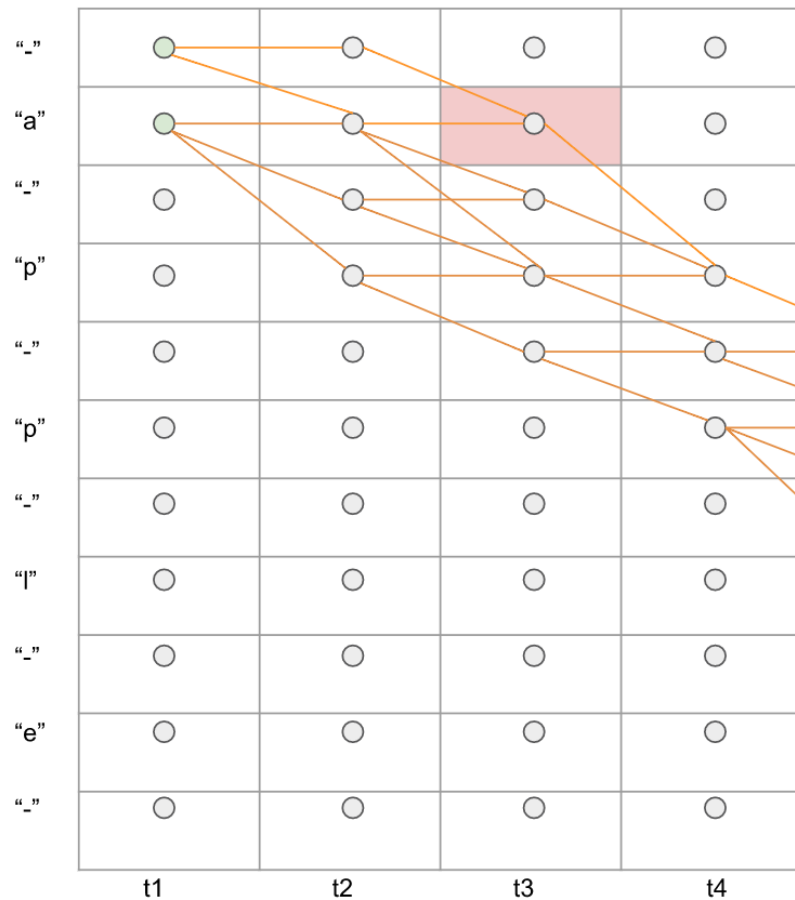
- Делается аналогично прямому проходу для каждого момента времени и метки рассчитываем вероятность

$$\beta_t(s) \stackrel{\text{def}}{=} \sum_{\substack{\pi \in N^T: \\ \mathcal{B}(\pi_{t:T}) = \mathbf{1}_{s:|I|}}} \prod_{t'=t}^T y_{\pi_{t'}}^{t'}$$

$$\begin{aligned} \beta_T(|I'|) &= y_b^T \\ \beta_T(|I'| - 1) &= y_{\mathbf{1}_{|I|}}^T \\ \beta_T(s) &= 0, \quad \forall s < |I'| - 1 \end{aligned}$$

$$\beta_t(s) = \begin{cases} \bar{\beta}_t(s) y_{I'_s}^t & \text{if } I'_s = b \text{ or } I'_{s+2} = I'_s \\ (\bar{\beta}_t(s) + \beta_{t+1}(s+2)) y_{I'_s}^t & \text{otherwise} \end{cases} \quad \bar{\beta}_t(s) \stackrel{\text{def}}{=} \beta_{t+1}(s) + \beta_{t+1}(s+1).$$

Вероятность путей для отдельной метки



$$\alpha_3(2) = p("--a") + p("-aa") + p("aaa") =$$

$$= y_-^1 \cdot y_-^2 \cdot y_a^3 + y_-^1 \cdot y_a^2 \cdot y_a^3 + y_a^1 \cdot y_a^2 \cdot y_a^3$$

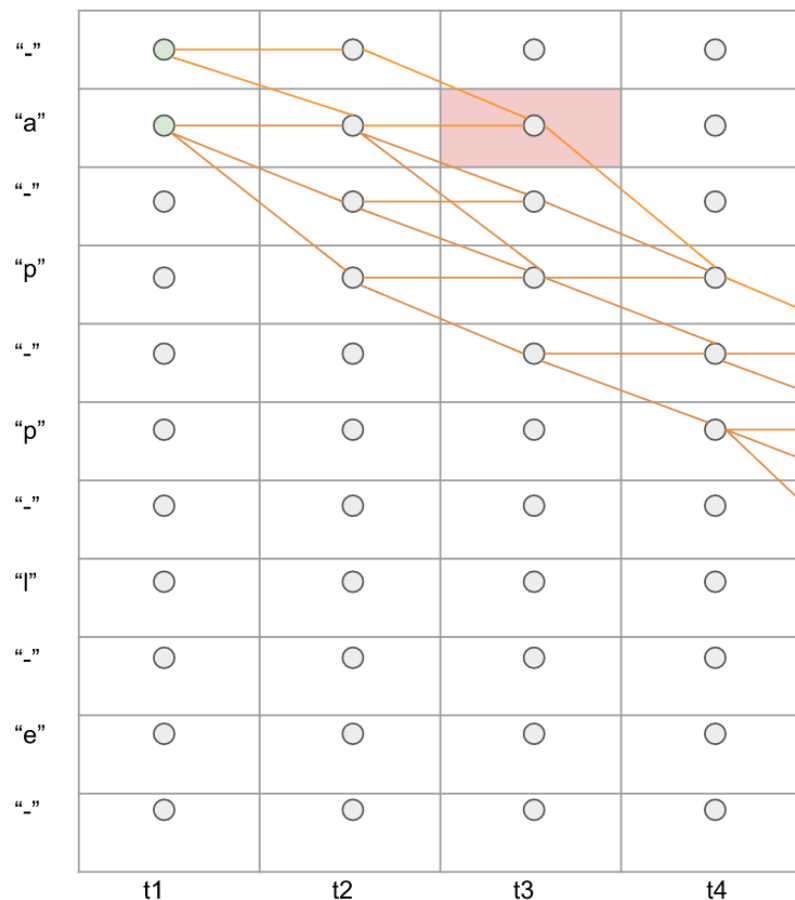
$$\beta_3(2) = p("ap-le") = y_a^3 \cdot y_p^4 \cdot y_-^5 \cdot y_p^6 \cdot y_l^7 \cdot y_e^8$$

$$\Rightarrow \alpha_3(2) \cdot \beta_3(2) = y_-^1 \cdot y_-^2 \cdot y_a^3 \cdot y_a^3 \cdot y_p^4 \cdot y_-^5 \cdot y_p^6 \cdot y_l^7 \cdot y_e^8 +$$

$$+ y_-^1 \cdot y_a^2 \cdot y_a^3 \cdot y_a^3 \cdot y_p^4 \cdot y_-^5 \cdot y_p^6 \cdot y_l^7 \cdot y_e^8$$

$$+ y_a^1 \cdot y_a^2 \cdot y_a^3 \cdot y_a^3 \cdot y_p^4 \cdot y_-^5 \cdot y_p^6 \cdot y_l^7 \cdot y_e^8$$

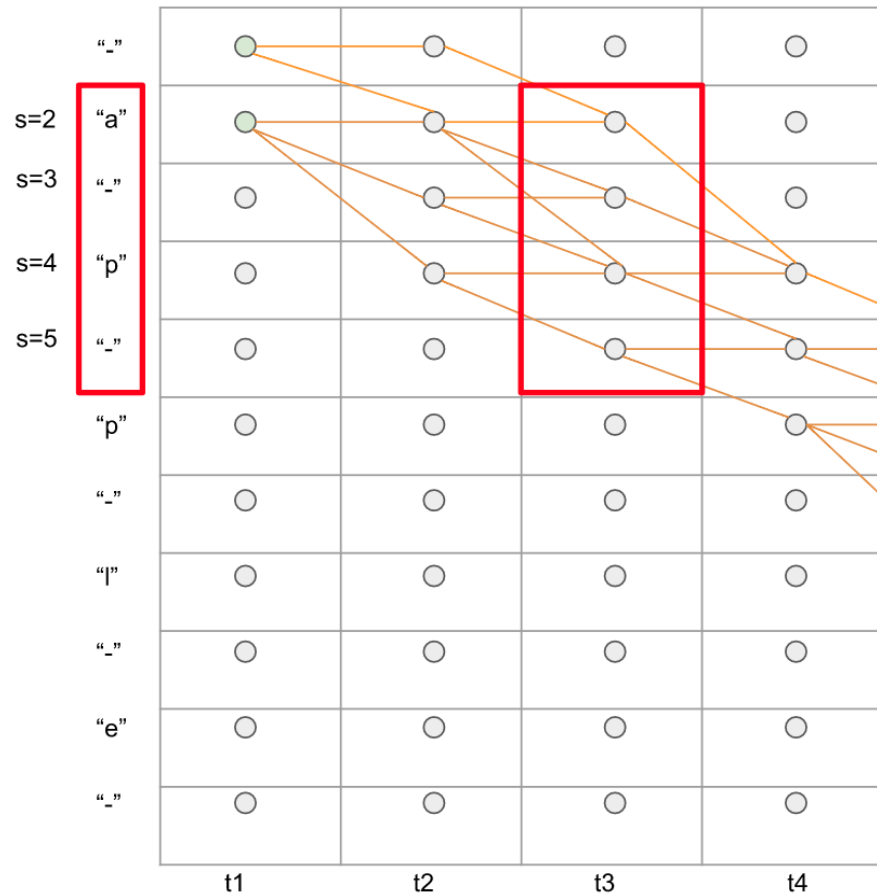
Вероятность путей для отдельной метки



$$\begin{aligned}
 \alpha_3(2) \cdot \beta_3(2) &= y_-^1 \cdot y_-^2 \cdot y_a^3 \cdot \boxed{y_a^3} \cdot y_p^4 \cdot y_-^5 \cdot y_p^6 \cdot y_l^7 \cdot y_e^8 + \\
 &+ y_-^1 \cdot y_a^2 \cdot y_a^3 \cdot \boxed{y_a^3} \cdot y_p^4 \cdot y_-^5 \cdot y_p^6 \cdot y_l^7 \cdot y_e^8 \\
 &+ y_a^1 \cdot y_a^2 \cdot y_a^3 \cdot \boxed{y_a^3} \cdot y_p^4 \cdot y_-^5 \cdot y_p^6 \cdot y_l^7 \cdot y_e^8 = \\
 &= (p("--aap-ple") + p("-aap-ple") + p("aap-ple")) * y_a^3
 \end{aligned}$$

$$\frac{\alpha_3(2) \cdot \beta_3(2)}{y_a^3} = (p("--aap-ple") + p("-aap-ple") + p("aap-ple"))$$

Вероятность слова для времени



Суммарная вероятность
всех путей на шаге 3

$$p(\text{"apple"}) = \sum_{s=2}^5 \frac{\alpha_3(s) \cdot \beta_3(s)}{y_{seq(s)}^3}$$

Вероятность в любой
момент времени

$$p(\text{"apple"}) = \sum_{s=1}^{|seq|} \frac{\alpha_t(s) \cdot \beta_t(s)}{y_{seq(s)}^t}$$

$$\text{CTC Loss} = -\ln(p(\text{"apple"}))$$

Градиент

$$p(\text{"apple"}) = \sum_{s=1}^{|\text{seq}|} \frac{\alpha_t(s) \cdot \beta_t(s)}{y_{\text{seq}(s)}^t}$$

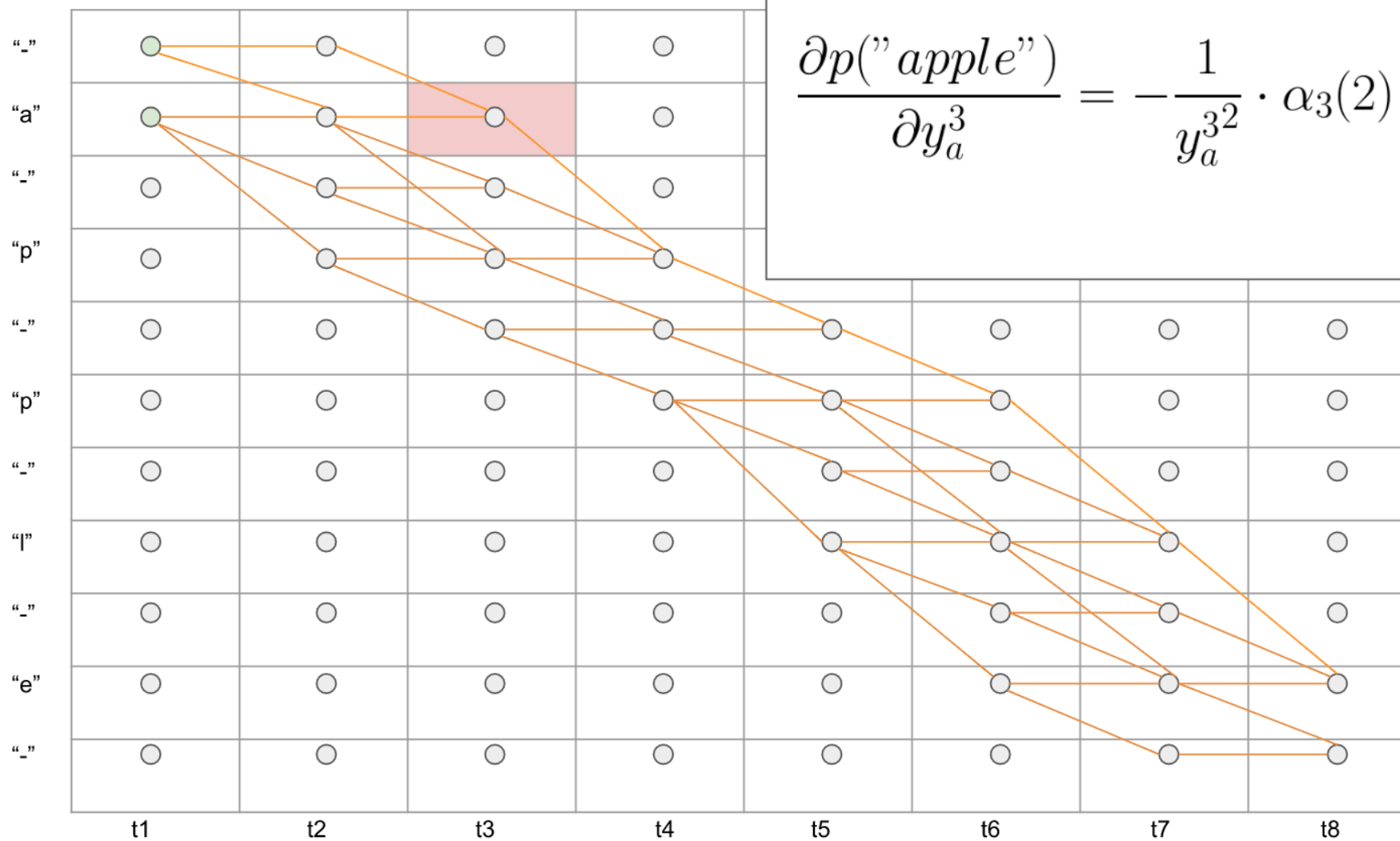
- Принимаем во внимание, что

$$\frac{\partial(-\ln(p(\text{"apple"})))}{\partial y_k^t} = -\frac{1}{p(\text{"apple"})} \cdot \frac{\partial p(\text{"apple"})}{\partial y_k^t}$$

- Чтобы обучить сеть нам нужно продифференцировать функцию ошибки по всем выходам сети y_k^t
- Рассматриваем только пути идущие через символ k в момент t

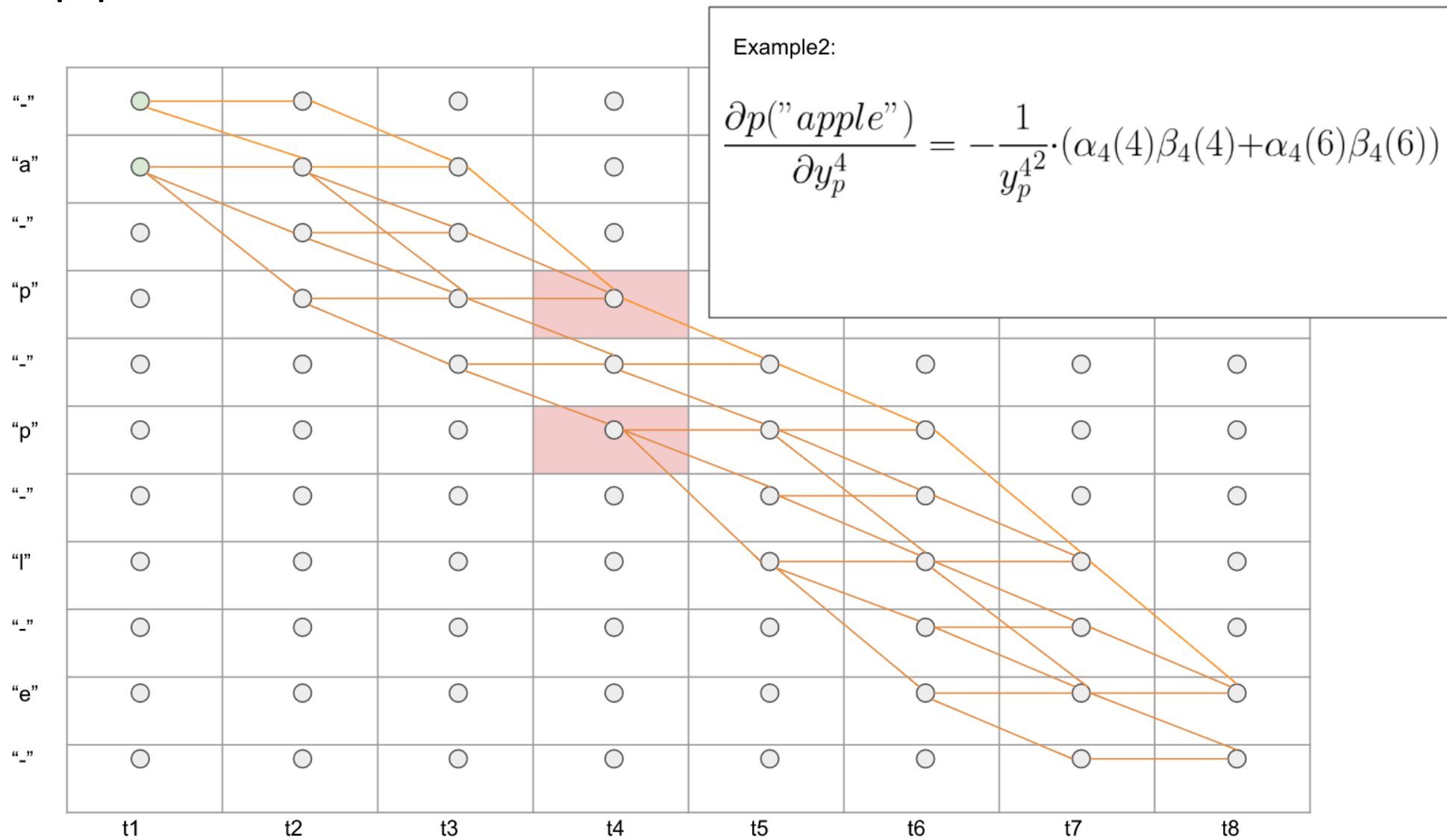
$$\frac{\partial p(\text{"apple"})}{\partial y_k^t} = -\frac{1}{y_k^{t2}} \cdot \sum_{s:\text{seq}(s)=k} \alpha_t(s) \cdot \beta_t(s)$$

Градиент



$$\frac{\partial p(\text{"apple"})}{\partial y_a^3} = -\frac{1}{y_a^{32}} \cdot \alpha_3(2) \cdot \beta_3(2)$$

Градиент



Deep Speech

- Вход – транскрибированное аудио самплы (частота дискретизации 8К)

$$\mathcal{X} = \tilde{\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots\}}.$$

- Фичи – спектрограммы. 80 линейно распределенных логарифмических фильтр-банков.

Deep Speech

- Архитектура сети

$$h_{t,k}^{(6)} = \hat{y}_{t,k} \equiv \mathbb{P}(c_t = k|x) = \frac{\exp(W_k^{(6)} h_t^{(5)} + b_k^{(6)})}{\sum_j \exp(W_j^{(6)} h_t^{(5)} + b_j^{(6)})}$$

Двунаправленная RNN

$$h_t^{(f)} = g(W^{(4)} h_t^{(3)} + W_r^{(f)} h_{t-1}^{(f)} + b^{(4)})$$

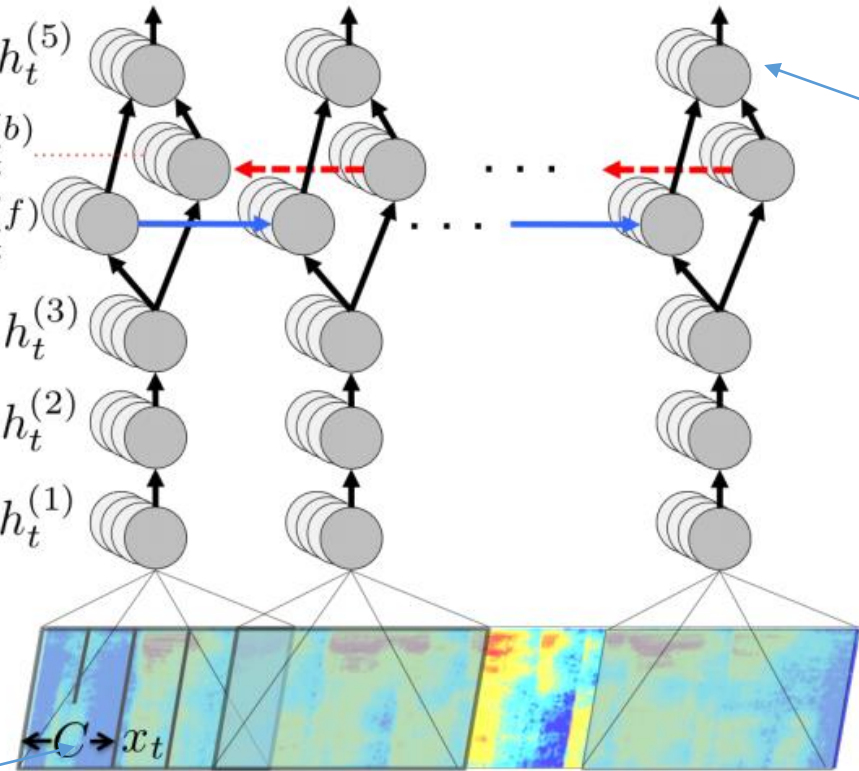
$$h_t^{(b)} = g(W^{(4)} h_t^{(3)} + W_r^{(b)} h_{t+1}^{(b)} + b^{(4)})$$

Линейный слой

$$h_t^{(l)} = g(W^{(l)} h_t^{(l-1)} + b^{(l)})$$

Линейный слой
 $g(W^{(5)} h_t^{(4)} + b^{(5)})$

$g = \text{ReLU}$



Количество сгруппированных фреймов
 $C \in \{5, 7, 9\}$

Deep Speech

- Предсказывает последовательность СИМВОЛОВ
- Loss = CTC
- Для предсказания использует языковую модель

$$Q(c) = \log(\mathbb{P}(c|x)) + \alpha \log(\mathbb{P}_{lm}(c)) + \beta \text{word_count}(c)$$

Вероятность
последователь
ности RNN

Вероятность последовательности языковой модели

RNN output	Decoded Transcription
what is the weather like in bostin right now prime miniter nerendr modi arther n tickets for the game	what is the weather like in boston right now prime minister narendra modi are there any tickets for the game