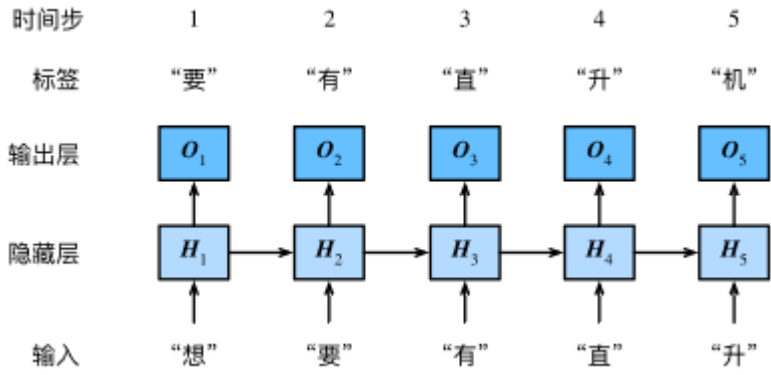


循环神经网络

本节介绍循环神经网络，下图展示了如何基于循环神经网络实现语言模型。我们的目的是基于当前的输入与过去的输入序列，预测序列的下一个字符。循环神经网络引入一个隐藏变量 \boldsymbol{H} ，用 \boldsymbol{H}_t 表示 \boldsymbol{H} 在时间步 t 的值。 \boldsymbol{H}_t 的计算基于 \boldsymbol{X}_t 和 \boldsymbol{H}_{t-1} ，可以认为 \boldsymbol{H}_t 记录了到当前字符为止的序列信息，利用 \boldsymbol{H}_t 对序列的下一个字符进行预测。



循环神经网络的构造

我们先看循环神经网络的具体构造。假设 $\boldsymbol{X}_t \in \mathbb{R}^{n \times d}$ 是时间步 t 的小批量输入， $\boldsymbol{H}_t \in \mathbb{R}^{n \times h}$ 是该时间步的隐藏变量，则：

$$\boldsymbol{H}_t = \phi(\boldsymbol{X}_t \boldsymbol{W}_{xh} + \boldsymbol{H}_{t-1} \boldsymbol{W}_{hh} + \boldsymbol{b}_h).$$

其中， $\boldsymbol{W}_{xh} \in \mathbb{R}^{d \times h}$ ， $\boldsymbol{W}_{hh} \in \mathbb{R}^{h \times h}$ ， $\boldsymbol{b}_h \in \mathbb{R}^{1 \times h}$ ， ϕ 函数是非线性激活函数。由于引入了 $\boldsymbol{H}_{t-1} \boldsymbol{W}_{hh}$ ， \boldsymbol{H}_t 能够捕捉截至当前时间步的序列的历史信息，就像是神经网络当前时间步的状态或记忆一样。由于 \boldsymbol{H}_t 的计算基于 \boldsymbol{H}_{t-1} ，上式的计算是循环的，使用循环计算的网路即循环神经网络（recurrent neural network）。

在时间步 t ，输出层的输出为：

$$\boldsymbol{O}_t = \boldsymbol{H}_t \boldsymbol{W}_{hq} + \boldsymbol{b}_q.$$

其中 $\boldsymbol{W}_{hq} \in \mathbb{R}^{h \times q}$ ， $\boldsymbol{b}_q \in \mathbb{R}^{1 \times q}$ 。

从零开始实现循环神经网络

我们先尝试从零开始实现一个基于字符级循环神经网络的语言模型，这里我们使用周杰伦的歌词作为语料，首先我们读入数据：

```
In [2]:
import torch
import torch.nn as nn
import time
import math
import sys
sys.path.append("/home/kesci/input")
import d2l_jay2900 as d2l
(corpus_indices, char_to_idx, idx_to_char, vocab_size) = d2l.load_data_jay_lyrics()
device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
```

one-hot向量

我们需要将字符表示成向量，这里采用one-hot向量。假设词典大小是 N ，每次字符对应一个从0到 $N - 1$ 的唯一的索引，则该字符的向量是一个长度为 N 的向量，若字符的索引是 i ，则该向量的第 i 个位置为1，其他位置为0。下面分别展示了索引为0和2的one-hot向量，向量长度等于词典大小。

```
In [3]:
def one_hot(x, n_class, dtype=torch.float32):
    result = torch.zeros(x.shape[0], n_class, dtype=dtype, device=x.device) # shape: (n, n_class)
    result.scatter_(1, x.long().view(-1, 1), 1) # result[i, x[i, 0]] = 1
    return result

x = torch.tensor([0, 2])
x_one_hot = one_hot(x, vocab_size)
print(x_one_hot)
print(x_one_hot.shape)
print(x_one_hot.sum(axis=1))

tensor([[1., 0., 0., ..., 0., 0., 0.],
        [0., 0., 1., ..., 0., 0., 0.]])
torch.Size([2, 1027])
tensor([1., 1.]])
```

我们每次采样的小批量的形状是（批量大小, 时间步数）。下面的函数将这样的小批量变换成数个形状为（批量大小, 词典大小）的矩阵，矩阵个数等于时间步数。也就是说，时间步 t 的输入为 $\boldsymbol{X}_t \in \mathbb{R}^{n \times d}$ ，其中 n 为批量大小， d 为词向量大小，即one-hot向量长度（词典大小）。

```
In [4]:
def to_onehot(X, n_class):
    return [one_hot(X[:, i], n_class) for i in range(X.shape[1])]

X = torch.arange(10).view(2, 5)
inputs = to_onehot(X, vocab_size)
print(len(inputs), inputs[0].shape)

5 torch.Size([2, 1027])
```

初始化模型参数

```
In [5]:
num_inputs, num_hidden, num_outputs = vocab_size, 256, vocab_size
# num_inputs: d
# num_hidden: h, 隐藏单元的个数是超参数
# num_outputs: q

def get_params():
    def _one(shape):
        param = torch.zeros(shape, device=device, dtype=torch.float32)
        nn.init.normal_(param, 0, 0.01)
        return torch.nn.Parameter(param)

    # 隐藏层参数
    W_xh = _one((num_inputs, num_hidden))
    W_hh = _one((num_hidden, num_hidden))
    b_h = torch.nn.Parameter(torch.zeros(num_hidden, device=device))
    # 输出层参数
    W_hq = _one((num_hidden, num_outputs))
    b_q = torch.nn.Parameter(torch.zeros(num_outputs, device=device))
    return (W_xh, W_hh, b_h, W_hq, b_q)
```

定义模型

函数rnn用循环的方式依次完成循环神经网络每个时间步的计算。

```
In [6]:
def rnn(inputs, state, params):
    # inputs和outputs皆为num_steps个形状为(batch_size, vocab_size)的矩阵
    W_xh, W_hh, b_h, W_hq, b_q = params
    H, = state
    outputs = []
    for X in inputs:
        H = torch.tanh(torch.matmul(X, W_xh) + torch.matmul(H, W_hh) + b_h)
        Y = torch.matmul(H, W_hq) + b_q
        outputs.append(Y)
    return outputs, (H,)
```

函数init_rnn_state初始化隐藏变量，这里的返回值是一个元组。

```
In [7]:
def init_rnn_state(batch_size, num_hidden, device):
    return (torch.zeros((batch_size, num_hidden), device=device), )
```

做个简单的测试来观察输出结果的个数（时间步数），以及第一个时间步的输出层输出的形状和隐藏状态的形状。

```
In [8]:
print(X.shape)
print(num_hidden)
print(vocab_size)
state = init_rnn_state(X.shape[0], num_hidden, device)
inputs = to_onehot(X.to(device), vocab_size)
params = get_params()
outputs, state_new = rnn(inputs, state, params)
print(len(inputs), inputs[0].shape)
print(len(outputs), outputs[0].shape)
print(len(state), state[0].shape)
print(len(state_new), state_new[0].shape)

torch.Size([2, 5])
256
1027
5 torch.Size([2, 1027])
5 torch.Size([2, 1027])
1 torch.Size([2, 256])
1 torch.Size([2, 256])
```

裁剪梯度

循环神经网络中较容易出现梯度衰减或梯度爆炸，这会导致网络几乎无法训练。裁剪梯度（clip gradient）是一种应对梯度爆炸的方法。假设我们把所有模型参数的梯度拼接成一个向量 \boldsymbol{g} ，并设裁剪的阈值是 θ 。裁剪后的梯度

$$\min\left(\frac{\theta}{\|\boldsymbol{g}\|}, 1\right)\boldsymbol{g}$$

的 L_2 范数不超过 θ 。

```
In [9]:
def grad_clipping(params, theta, device):
    norm = torch.tensor([0.0], device=device)
    for param in params:
        norm += (param.grad.data ** 2).sum()
    norm = norm.sqrt().item()
    if norm > theta:
        for param in params:
            param.grad.data *= (theta / norm)
```

定义预测函数

以下函数基于前缀prefix（含有数个字符的字符串）来预测接下来的num_chars个字符。这个函数稍显复杂，其中我们将循环神经单元rnn设置成了函数参数，这样在后面小节介绍其他循环神经网络时能重复使用这个函数。

In [10]:

```
def predict_rnn(prefix, num_chars, rnn, params, init_rnn_state,
                num_hiddens, vocab_size, device, idx_to_char, char_to_idx):
    state = init_rnn_state(1, num_hiddens, device)
    output = [char_to_idx[prefix[0]]] # output记录prefix加上预测的num_chars个字符
    for t in range(num_chars + len(prefix) - 1):
        # 将上一时间步的输出作为当前时间步的输入
        X = to_onehot(torch.tensor([[output[-1]]], device=device), vocab_size)
        # 计算输出和更新隐藏状态
        (Y, state) = rnn(X, state, params)
        # 下一个时间步的输入是prefix里的字符或者当前的最佳预测字符
        if t < len(prefix) - 1:
            output.append(char_to_idx[prefix[t + 1]])
        else:
            output.append(Y[0].argmax(dim=1).item())
    return ''.join([idx_to_char[i] for i in output])
```

我们先测试一下predict_rnn函数。我们将根据前缀“分开”创作长度为10个字符（不考虑前缀长度）的一段歌词。因为模型参数为随机值，所以预测结果也是随机的。

In [11]:

```
predict_rnn('分开', 10, rnn, params, init_rnn_state, num_hiddens, vocab_size,
            device, idx_to_char, char_to_idx)
```

Out[11]:

'分开却难蛎知化落果够辛抢'

困惑度

我们通常使用困惑度（perplexity）来评价语言模型的好坏。回忆一下“[softmax回归](#)”一节中交叉熵损失函数的定义。困惑度是对交叉熵损失函数做指数运算后得到的值。特别地，

- 最佳情况下，模型总是把标签类别的概率预测为1，此时困惑度为1；
- 最坏情况下，模型总是把标签类别的概率预测为0，此时困惑度为正无穷；
- 基线情况下，模型总是预测所有类别的概率都相同，此时困惑度为类别个数。

显然，任何一个有效模型的困惑度必须小于类别个数。在本例中，困惑度必须小于词典大小vocab_size。

定义模型训练函数

跟之前章节的模型训练函数相比，这里的模型训练函数有以下几点不同：

1. 使用困惑度评价模型。
2. 在迭代模型参数前裁剪梯度。
3. 对时序数据采用不同采样方法将导致隐藏状态初始化的不同。

In [12]:

```
def train_and_predict_rnn(rnn, get_params, init_rnn_state, num_hiddens,
                          vocab_size, device, corpus_indices, idx_to_char,
                          char_to_idx, is_random_iter, num_epochs, num_steps,
                          lr, clipping_theta, batch_size, pred_period,
                          pred_len, prefixes):

    if is_random_iter:
        data_iter_fn = d2l.data_iter_random
    else:
        data_iter_fn = d2l.data_iter_consecutive
    params = get_params()
    loss = nn.CrossEntropyLoss()

    for epoch in range(num_epochs):
        if not is_random_iter: # 如使用相邻采样，在epoch开始时初始化隐藏状态
            state = init_rnn_state(batch_size, num_hiddens, device)
        l_sum, n, start = 0.0, 0, time.time()
        data_iter = data_iter_fn(corpus_indices, batch_size, num_steps, device)
        for X, Y in data_iter:
            if is_random_iter: # 如使用随机采样，在每个小批量更新前初始化隐藏状态
                state = init_rnn_state(batch_size, num_hiddens, device)
            else: # 否则需要使用detach函数从计算图分离隐藏状态
                for s in state:
                    s.detach_()

            # inputs是num_steps个形状为(batch_size, vocab_size)的矩阵
            inputs = to_onehot(X, vocab_size)
            # outputs有num_steps个形状为(batch_size, vocab_size)的矩阵
            (outputs, state) = rnn(inputs, state, params)
            # 拼接之后形状为(num_steps * batch_size, vocab_size)
            outputs = torch.cat(outputs, dim=0)
            # Y的形状是(batch_size, num_steps)，转置后再变成形状为
            # (num_steps * batch_size,)的向量，这样跟输出的行一一对应
            y = torch.flatten(Y.T)
            # 使用交叉熵损失计算平均分类误差
            l = loss(outputs, y.long())

            # 梯度清0
            if params[0].grad is not None:
                for param in params:
                    param.grad.data.zero_()
            l.backward()
            grad_clipping(params, clipping_theta, device) # 裁剪梯度
            d2l.sgd(params, lr, 1) # 因为误差已经取过均值，梯度不用再做平均
            l_sum += l.item() * y.shape[0]
            n += y.shape[0]

        if (epoch + 1) % pred_period == 0:
            print('epoch %d, perplexity %f, time %.2f sec' % (
                epoch + 1, math.exp(l_sum / n), time.time() - start))
            for prefix in prefixes:
                print('-', predict_rnn(prefix, pred_len, rnn, params, init_rnn_state,
                    num_hiddens, vocab_size, device, idx_to_char, char_to_idx))
```

训练模型并创作歌词

现在我们可以训练模型了。首先，设置模型超参数。我们将根据前缀“分开”和“不分开”分别创作长度为50个字符（不考虑前缀长度）的一段歌词。我们每过50个迭代周期便根据当前训练的模型创作一段歌词。

In [13]:

```
num_epochs, num_steps, batch_size, lr, clipping_theta = 250, 35, 32, 1e-2
pred_period, pred_len, prefixes = 50, 50, ['分开', '不分开']
```

下面采用随机采样训练模型并创作歌词。

In [14]:

```
train_and_predict_rnn(rnn, get_params, init_rnn_state, num_hiddens,
                      vocab_size, device, corpus_indices, idx_to_char,
                      char_to_idx, True, num_epochs, num_steps, lr,
                      clipping_theta, batch_size, pred_period, pred_len,
                      prefixes)

epoch 50, perplexity 69.143026, time 0.59 sec
- 分开 我不要再想 我不要再想 我不要再想 我不要再想 我不要再想 我不要再想 我不要再想 我
- 不分开 我有那 我想么这 我不要再 你的让我 你的让我 你的让我 你的让我 你的让我 你的让我
epoch 100, perplexity 9.620708, time 0.50 sec
- 分开 一颗两双 在谁心空 我不定再不 没有你的我有 思不要 后情走 我想就这样牵着你的手不放开 爱可不
- 不分开堡 我爱你这你 我不 你不 我不能再想 我不 再不 我不 我不 我不 我不 我不 我不 我不
epoch 150, perplexity 2.758654, time 0.53 sec
- 分开 一只到 一颗两步三步四步望著天 看星星 一颗两颗三颗四颗 连成线背著背默默许下心愿 看远方的星是否
- 不分开扫 我不能再想 我不 我不 我不要 爱情走的太快就像龙卷风 不能承受我已无处可躲 我不要再想 我不能
epoch 200, perplexity 1.563496, time 0.57 sec
- 分开 一只到它心手的母斑鸠 在的黑猫笑走 学你去没有 静单三 教给堂 什么四对 在漠在中怎么擅有的窝
- 不分开想 我不能再想 我不 我不 我不能 爱情走的太快就像龙卷风 不能承受我已无处可躲 我不要再想 我不要
epoch 250, perplexity 1.304408, time 0.50 sec
- 分开 一直到 因手堂里术光 折变黏的豆都 然梅成著我 一口一口吃掉忧愁 别些了午三点阳光 瞎教太危重泪
- 不分开吗 我叫你爸 你打我妈 这样对吗干嘛这样 何必让酒牵鼻子走 瞎 说一个痛废 就这伦现不 穿不想烦
```

接下来采用相邻采样训练模型并创作歌词。

In [15]:

```
train_and_predict_rnn(rnn, get_params, init_rnn_state, num_hiddens,
                      vocab_size, device, corpus_indices, idx_to_char,
                      char_to_idx, False, num_epochs, num_steps, lr,
                      clipping_theta, batch_size, pred_period, pred_len,
                      prefixes)

epoch 50, perplexity 59.785954, time 0.59 sec
- 分开 我不要的可爱女人 坏坏的让我疯狂的可爱女人 坏坏的让我疯狂的可爱女人 坏坏的让我疯狂的可爱女人 坏
- 不分开 我不要的可爱女人 坏成的让我疯狂的可爱女人 坏坏的让我疯狂的可爱女人 坏坏的让我疯狂的可爱女人 坏
epoch 100, perplexity 6.826339, time 0.50 sec
- 分开 我说好这你的我像 你为样口 整小村外的溪边 默默激待 一颗 有什么不 再来一碗热粥 配上几客栈人多
- 不分开觉 我已经再你看 这生 却又再考倒吧 想通 你想很久了吧 通 却对很久 戒指在 分片怎么 快些有
epoch 150, perplexity 2.071869, time 0.50 sec
- 分开 我说儿这你怎到你的手不放开 爱可不可以简简单单没有伤害 你 靠着我的肩膀 你 在我胸口睡著 像这样
- 不分开觉透你不经 说 让不多口听是我 说不你 语沉默 一子四三 每头一碗的粥边 默默等待 娘子 有什么不妥
epoch 200, perplexity 1.290922, time 0.42 sec
- 分开 我说儿有你每了 却少林童武故事 有是堂有城堡 每天忙碌地的寻找 到底什么我想要 却发现迷了路怎么
- 不分开觉 你已经离开我 不知不觉 我跟了这节奏 后知后觉 又过了一个秋 后知后觉 我该好好生活 我该好好生
epoch 250, perplexity 1.150748, time 0.53 sec
- 分开 我说好这你怎么 却少林跟武当 快使用双截棍 哼哼哈兮 快使用双截棍 哼哼哈兮 如果我有轻功 飞檐走
- 不分开觉 你已经离开我 不知不觉 我跟了这节奏 后知后觉 后知了觉 迷迷蒙蒙 你给的梦 出现裂缝 隐隐作痛
```

循环神经网络的简介实现

定义模型

我们使用Pytorch中的nn.RNN来构造循环神经网络。在本节中，我们主要关注nn.RNN的以下几个构造函数参数：

- input_size - The number of expected features in the input x
- hidden_size – The number of features in the hidden state h
- nonlinearity – The non-linearity to use. Can be either 'tanh' or 'relu'. Default: 'tanh'
- batch_first – If True, then the input and output tensors are provided as (batch_size, num_steps, input_size). Default: False

这里的batch_first决定了输入的形状，我们使用默认的参数False，对应的输入形状是 (num_steps, batch_size, input_size)。

forward函数的参数为：

- input of shape (num_steps, batch_size, input_size): tensor containing the features of the input sequence.
- h_0 of shape (num_layers * num_directions, batch_size, hidden_size): tensor containing the initial hidden state for each element in the batch. Defaults to zero if not provided. If the RNN is bidirectional, num_directions should be 2, else it should be 1.

forward函数的返回值是：

- output of shape (num_steps, batch_size, num_directions * hidden_size): tensor containing the output features (h_t) from the last layer of the RNN, for each t.
- h_n of shape (num_layers * num_directions, batch_size, hidden_size): tensor containing the hidden state for t = num_steps.

现在我们构造一个nn.RNN实例，并用一个简单的例子来看一下输出的形状。

In [16]:

```
rnn_layer = nn.RNN(input_size=vocab_size, hidden_size=num_hiddens)
num_steps, batch_size = 35, 2
X = torch.rand(num_steps, batch_size, vocab_size)
state = None
Y, state_new = rnn_layer(X, state)
print(Y.shape, state_new.shape)

torch.Size([35, 2, 256]) torch.Size([1, 2, 256])
```

我们定义一个完整的基于循环神经网络的语言模型。

In [17]:

```
class RNNModel(nn.Module):
    def __init__(self, rnn_layer, vocab_size):
        super(RNNModel, self).__init__()
        self.rnn = rnn_layer
        self.hidden_size = rnn_layer.hidden_size * (2 if rnn_layer.bidirectional else 1)
        self.vocab_size = vocab_size
        self.dense = nn.Linear(self.hidden_size, vocab_size)

    def forward(self, inputs, state):
        # inputs.shape: (batch_size, num_steps)
        X = to_onehot(inputs, vocab_size)
        X = torch.stack(X) # X.shape: (num_steps, batch_size, vocab_size)
        hiddens, state = self.rnn(X, state)
        hiddens = hiddens.view(-1, hiddens.shape[-1]) # hiddens.shape: (num_steps * batch_size, hidden_size)
        output = self.dense(hiddens)
        return output, state
```

类似的，我们需要实现一个预测函数，与前面的区别在于前向计算和初始化隐藏状态。

In [18]:

```
def predict_rnn_pytorch(prefix, num_chars, model, vocab_size, device, idx_to_char,
                        char_to_idx):
    state = None
    output = [char_to_idx[prefix[0]]] # output记录prefix加上预测的num_chars个字符
    for t in range(num_chars + len(prefix) - 1):
        X = torch.tensor([output[-1]], device=device).view(1, 1)
        (Y, state) = model(X, state) # 前向计算不需要传入模型参数
        if t < len(prefix) - 1:
            output.append(char_to_idx[prefix[t + 1]])
        else:
            output.append(Y.argmax(dim=1).item())
    return ''.join([idx_to_char[i] for i in output])
```

使用权重为随机值的模型来预测一次。

In [19]:

```
model = RNNModel(rnn_layer, vocab_size).to(device)
predict_rnn_pytorch('分开', 10, model, vocab_size, device, idx_to_char, char_to_idx)
```

Out[19]:

'分开软火喝喝喝福喝福喝福'

接下来实现训练函数，这里只使用了相邻采样。

In [20]:

```
def train_and_predict_rnn_pytorch(model, num_hiddens, vocab_size, device,
                                corpus_indices, idx_to_char, char_to_idx,
                                num_epochs, num_steps, lr, clipping_theta,
                                batch_size, pred_period, pred_len, prefixes):
    loss = nn.CrossEntropyLoss()
    optimizer = torch.optim.Adam(model.parameters(), lr=lr)
    model.to(device)
    for epoch in range(num_epochs):
        l_sum, n, start = 0.0, 0, time.time()
        data_iter = d2l.data_iter_consecutive(corpus_indices, batch_size, num_steps, device) # 相邻采样
        state = None
        for X, Y in data_iter:
            if state is not None:
                # 使用detach函数从计算图分离隐藏状态
                if isinstance(state, tuple): # LSTM, state:(h, c)
                    state[0].detach_()
                    state[1].detach_()
                else:
                    state.detach_()
            (output, state) = model(X, state) # output.shape: (num_steps * batch_size, vocab_size)
            y = torch.flatten(Y.T)
            l = loss(output, y.long())

            optimizer.zero_grad()
            l.backward()
            grad_clipping(model.parameters(), clipping_theta, device)
            optimizer.step()
            l_sum += l.item() * y.shape[0]
            n += y.shape[0]

        if (epoch + 1) % pred_period == 0:
            print('epoch %d, perplexity %f, time %.2f sec' % (
                epoch + 1, math.exp(l_sum / n), time.time() - start))
            for prefix in prefixes:
                print(' -', predict_rnn_pytorch(
                    prefix, pred_len, model, vocab_size, device, idx_to_char,
                    char_to_idx))
```

训练模型。

In [21]:

```
num_epochs, batch_size, lr, clipping_theta = 250, 32, 1e-3, 1e-2
pred_period, pred_len, prefixes = 50, 50, ['分开', '不分开']
train_and_predict_rnn_pytorch(model, num_hiddens, vocab_size, device,
                               corpus_indices, idx_to_char, char_to_idx,
                               num_epochs, num_steps, lr, clipping_theta,
                               batch_size, pred_period, pred_len, prefixes)

epoch 50, perplexity 11.845853, time 0.41 sec
- 分开始想要 不能 你 我有了口让我感动的可爱女人 坏坏的让我疯狂的可爱女人 坏坏的让我疯狂的可爱女人 坏
- 不分开不 你不想我你 我不要再想 我不要再想 我不要再想 我不要再想 我不要再想 我不要再想 我不要再想
epoch 100, perplexity 1.272803, time 0.43 sec
- 分开不会像 风 没有伤害 你 靠着我的肩膀 你 在我胸口睡著 像这样的生活 我爱你 你爱我 我想大声宣
- 不分开不 像跟了很 泪真的牛 一跟好着 不慢的回忆 我有你看想 想要你却只剩大回忆 相爱还有别离 像无法
epoch 150, perplexity 1.065754, time 0.43 sec
- 分开不会像 也没有 烦我不有多烦恼 没有你烦我 有多烦恼 没有你烦 我有多难熬 没有你烦我有多
- 不分开不 像跟了这样的生活 爱你 我爱你 我想大声宣妈 对你依依不舍 连隔壁邻居都猜到我现在的感受 河边
epoch 200, perplexity 1.041097, time 0.40 sec
- 分开不会开 今我的别过 不透 我想我不想 说 就怎么每天祈就我的心跳你 不能 不知太觉 你已经离开我
- 不分开不 像跟了很 心伤的快动 它爱里 想一定一步云掉变 我面红默铁等待我季来临变沼泽 灰狼啃食著水鹿的
epoch 250, perplexity 1.021259, time 0.43 sec
- 分开不会开 了口我的成透 不开 你笑得 不懂 你已经很久 说你心在个人 你怎么每在人 什么心多 的雨
- 不分开不 像跟了这样的甜蜜 我在起痛里你的老斑 印地安老斑鸠 腿短毛不多 几天都没有喝水也能活 脑袋瓜
```