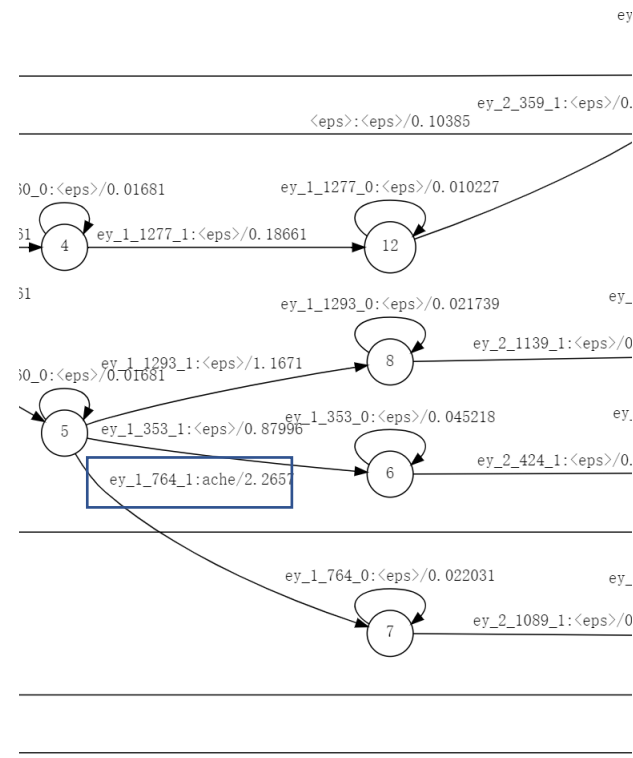
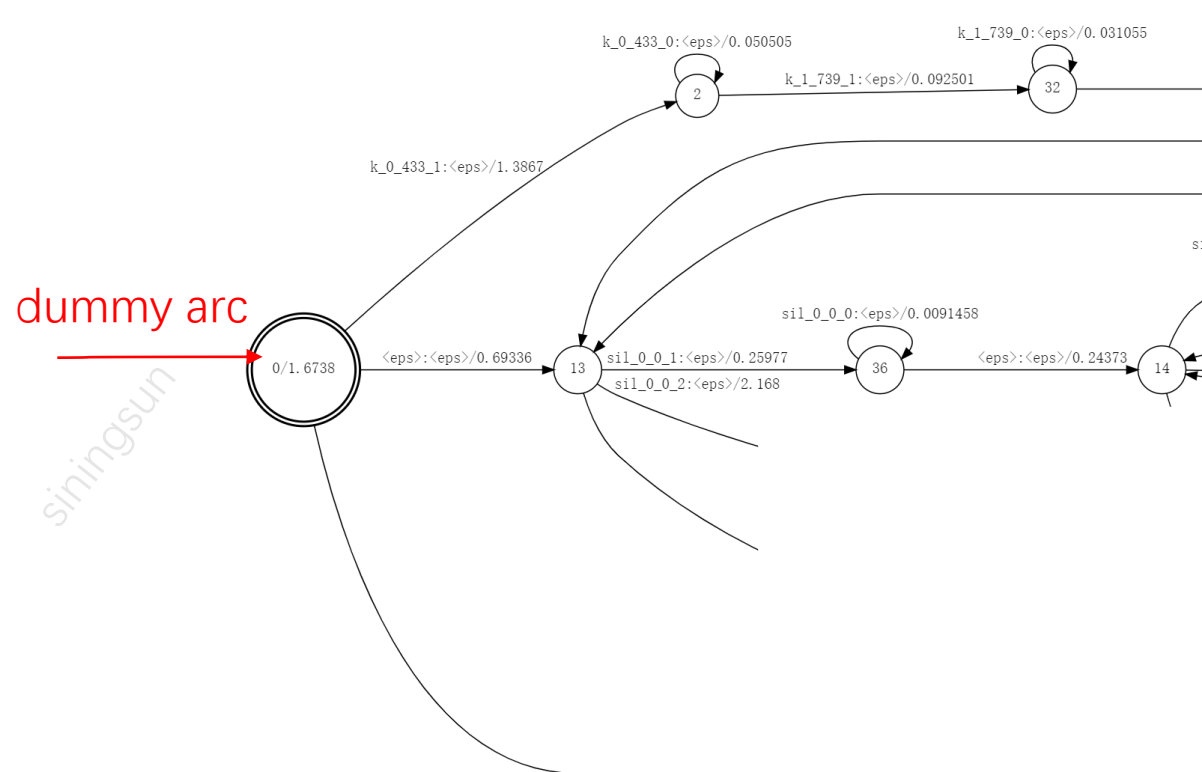


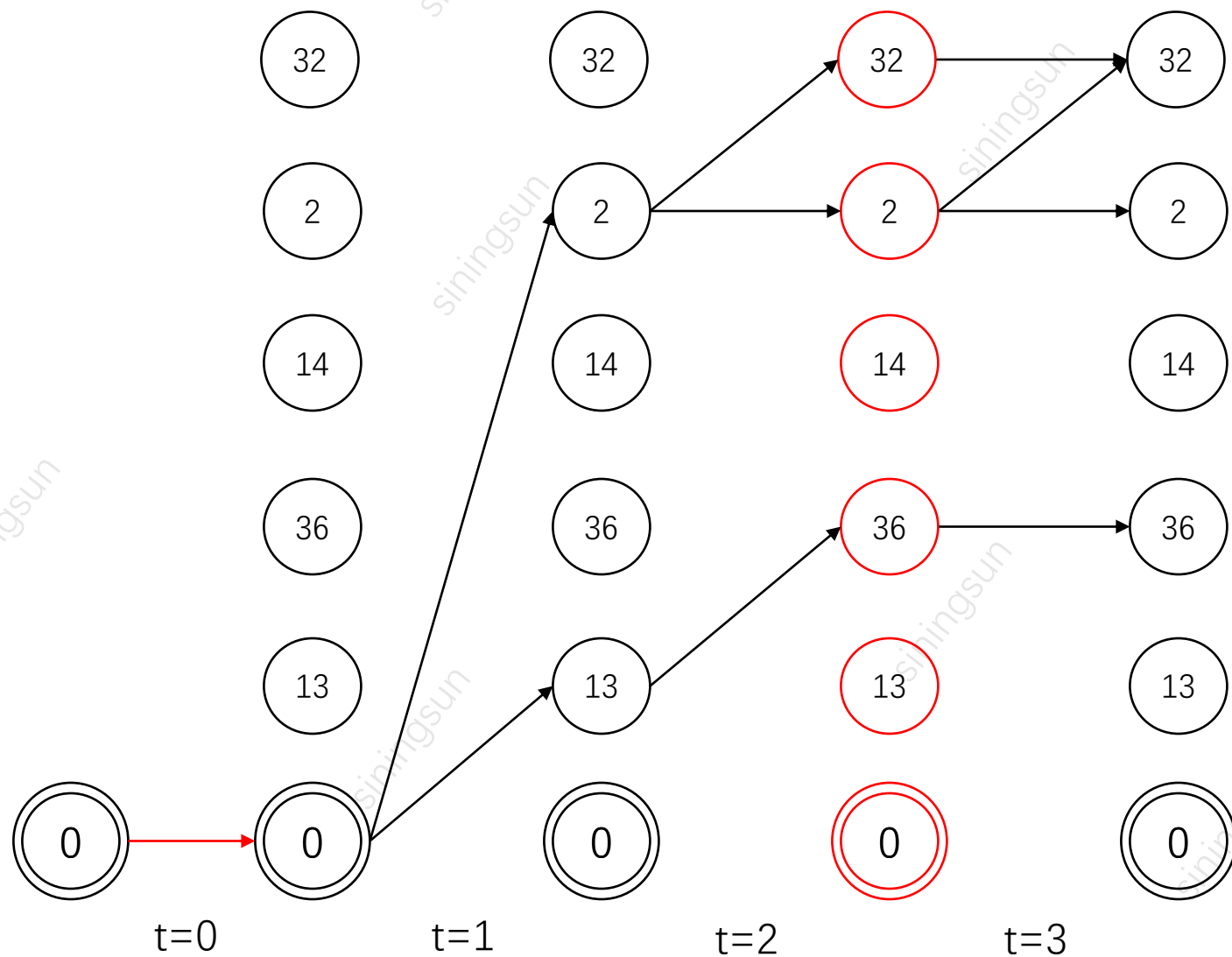
Kaldi SimpleDecoder 代码解析+逐步实例

孙思宁

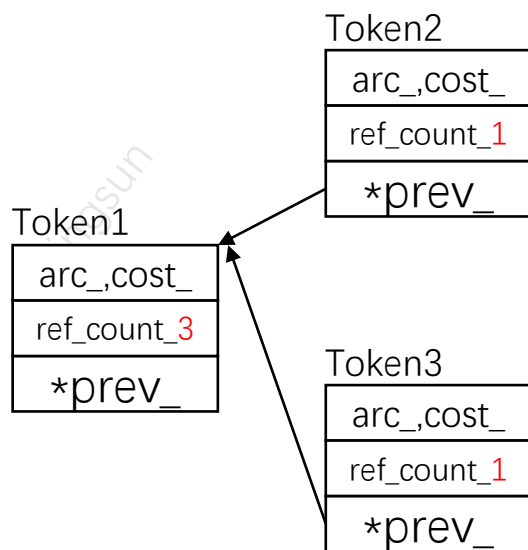
示例HCLG解码图



将HCLG沿着时间轴展开



arc_	HCLG这个WFST上的转移弧，弧上会有声学分、语言分、输入和输出label，以及目的状态(104-107)
*prev_	指向上一个Token的回溯指针，表示当前Token是从哪个Token传递过来的
ref_count_	主要是用来垃圾回收的，Token每当被其他Token引用一次，计数器+1(109)
cost_	cost_: 当前Token中累积分数
TokenDelete	只有ref_count_为1的Token才能被删除，说明它不会被其他Token回溯到；TokenDelete是递归回溯删除ref_count_等于1的整条路径上的Token



Token示例

Token结构体

InitDecoding

- 47行，cur_toks加入第一个token，对应图中Tokne1
- 第49行，处理从state0转移且输入label为空的所有的弧

```

93  class Token {
94  public:
95      LatticeArc arc_; // We use LatticeArc so that we can separately
96                      // store the acoustic and graph cost, in case
97                      // we need to produce lattice-formatted output.
98
99      Token *prev_;
100      int32 ref_count_;
101      double cost_; // accumulated total cost up to this point.
102
103      Token(const StdArc &arc,
104            BaseFloat acoustic_cost,
105            Token *prev): prev_(prev), ref_count_(1) {
106          arc_.ilabel = arc.ilabel;
107          arc_.olabel = arc.olabel;
108          arc_.weight = LatticeWeight(arc.weight.Value(), acoustic_cost);
109          arc_.nextstate = arc.nextstate;
110          if (prev) {
111              prev->ref_count_++;
112              cost_ = prev->cost_ + (arc.weight.Value() + acoustic_cost);
113          } else {
114              cost_ = arc.weight.Value() + acoustic_cost;
115          }
116      }
117
118      bool operator < (const Token &other) {
119          return cost_ > other.cost_;
120      }
121
122      static void TokenDelete(Token *tok) {
123          while (--tok->ref_count_ == 0) {
124              Token *prev = tok->prev_;
125              delete tok;
126              if (prev == NULL) return;
127              else tok = prev;
128          }
129      }
130  };
131
132  #ifndef KALDI_PARANOID
133      KALDI_ASSERT(tok->ref_count_ > 0);
134  #endif

```

prev_toks(空)

cur_toks_

32

2

36

13

0

Token1

arc_,cost_

ref_count_1

*prev_

t=0

t=1

t=2

逐步分解SimpleDecoder的解码过程

InitDecoding

1. 47行, cur_toks加入第一个token, 对应图中Tokne1
2. 第49行, 处理从state0转移且输入label为空的所有的弧

```
39 void SimpleDecoder::InitDecoding() {
40     // clean up from last time:
41     ClearToks(cur_toks_);
42     ClearToks(prev_toks_);
43     // initialize decoding:
44     StateId start_state = fst_.Start();
45     KALDI_ASSERT(start_state != fst::kNoStateId);
46     StdArc dummy_arc(0, 0, StdWeight::One(), start_state);
47     cur_toks_[start_state] = new Token(dummy_arc, 0.0, NULL);
48     num_frames_decoded_ = 0;
49     ProcessNonemitting();
50 }
51
```

prev_toks(空)

cur_toks_

32

2

36

13

0

Token2

arc_cost_

ref_count_1

*prev_

Token1

arc_cost_

ref_count_2

*prev_

← +1

t=0

t=1

t=2

逐步分解SimpleDecoder的解码过程

ProcessNonemitting()

1. 本例子只有0->13这一个满足要求的转移
2. 新建一个状态13的Token2，加入到cur_toks_
3. Token2指向了Token1，Token1的ref_count_+1

```
208 void SimpleDecoder::ProcessNonemitting() {
209     // Processes nonemitting arcs for one frame. Propagates within
210     // cur_toks_.
211     std::vector<StateId> queue;
212     double infinity = std::numeric_limits<double>::infinity();
213     double best_cost = infinity;
214     for (unordered_map<StateId, Token*>::iterator iter = cur_toks_.begin();
215          iter != cur_toks_.end();
216          ++iter) {
217         queue.push_back(iter->first);
218         best_cost = std::min(best_cost, iter->second->cost_);
219     }
220     double cutoff = best_cost + beam_;
221
222     while (!queue.empty()) {
223         StateId state = queue.back();
224         queue.pop_back();
225         Token *tok = cur_toks_[state];
226         KALDI_ASSERT(tok != NULL && state == tok->arc.nextstate);
227         for (fst::ArcIterator<fst::Fst<StdArc> > aiter(fst_, state);
228              !aiter.Done();
229              aiter.Next()) {
230             const StdArc &arc = aiter.Value();
231             if (arc.ilabel == 0) { // propagate nonemitting only.只处理ilabel为空的弧，空转移没有AM分数
232                 const BaseFloat acoustic_cost = 0.0;
233                 Token *new_tok = new Token(arc, acoustic_cost, tok);
234                 if (new_tok->cost_ > cutoff) {
235                     Token::TokenDelete(new_tok);
236                 } else {
237                     unordered_map<StateId, Token*>::iterator find_iter
238                         = cur_toks_.find(arc.nextstate);
239                     if (find_iter == cur_toks_.end()) {
240                         cur_toks_[arc.nextstate] = new_tok;
241                         queue.push_back(arc.nextstate);
242                     } else {
243                         if ( *(find_iter->second) < *new_tok ) {
244                             Token::TokenDelete(find_iter->second);
245                             find_iter->second = new_tok;
246                             queue.push_back(arc.nextstate);
247                         } else {
248                             Token::TokenDelete(new_tok);
249                         }
250                     }
251                 }
252             }
253         }
254     }
```

cur_toks_ (空)

prev_toks

32

2

36

13

0

Token2

arc_cost_
ref_count_1
*prev_

Token1

arc_cost_
ref_count_2
*prev_

t=0

t=1

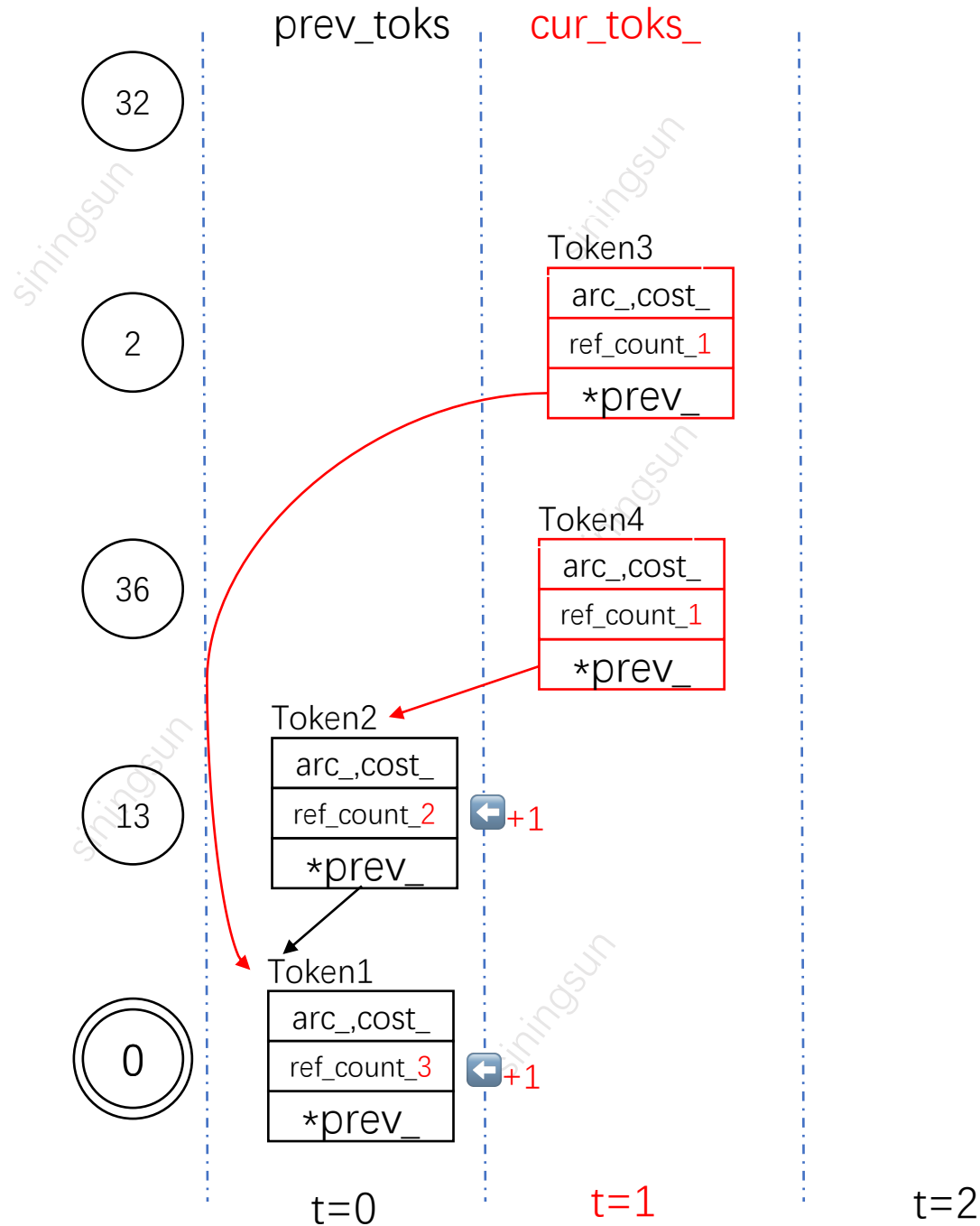
t=2

逐步分解SimpleDecoder的解码过程

AdvanceDecoding

1. 解码器开始逐帧处理decodable中的后验概率
2. while循环中，首先ClearToks，此时prev_toks_本身就是空；
3. 将cur_toks_和prev_toks_交换，然后70行开始核心代码

```
52 void SimpleDecoder::AdvanceDecoding(DecodableInterface *decodable,
53                                     int32 max_num_frames) {
54     KALDI_ASSERT(num_frames_decoded_ >= 0 &&
55                 "You must call InitDecoding() before AdvanceDecoding()");
56     int32 num_frames_ready = decodable->NumFramesReady();
57     // num_frames_ready must be >= num_frames_decoded, or else
58     // the number of frames ready must have decreased (which doesn't
59     // make sense) or the decodable object changed between calls
60     // (which isn't allowed).
61     KALDI_ASSERT(num_frames_ready >= num_frames_decoded_);
62     int32 target_frames_decoded = num_frames_ready;
63     if (max_num_frames >= 0)
64         target_frames_decoded = std::min(target_frames_decoded,
65                                         num_frames_decoded_ + max_num_frames);
66     while (num_frames_decoded_ < target_frames_decoded) {
67         // note: ProcessEmitting() increments num frames decoded
68         ClearToks(prev_toks_);
69         cur_toks_.swap(prev_toks_);
70         ProcessEmitting(decodable);
71         ProcessNonemitting();
72         PruneToks(beam_, &cur_toks_);
73     }
74 }
75
```



逐步分解SimpleDecoder的解码过程

ProcessEmitting()

- line 172-204, 遍历prev_toks_, 而从178-203, 遍历每个prev_toks_中的所有可以到达的非空转移
- 本例子中, S13->S36, S0->S2两个转移产生新的Token3和4, 放入了cur_toks, 注意此时Token1, 2的计数器都增加了1

```

167 void SimpleDecoder::ProcessEmitting(DecodableInterface *decodable) {
168     int32 frame = num_frames_decoded_;
169     // Processes emitting arcs for one frame. Propagates from
170     // prev_toks_ to cur_toks_.
171     double cutoff = std::numeric_limits<BaseFloat>::infinity();
172     for (unordered_map<StateId, Token*>::iterator iter = prev_toks_.begin();
173          iter != prev_toks_.end(); ++iter) {
174         ++iter) {
175             StateId state = iter->first;
176             Token *tok = iter->second;
177             KALDI_ASSERT(state == tok->arc_.nextstate);
178             for (fst::ArcIterator<fst::Fst<StdArc> > aiter(fst_, state);
179                  !aiter.Done(); aiter.Next()) {
180                 const StdArc &arc = aiter.Value();
181                 if (arc.ilabel != 0) { // propagate..
182                     BaseFloat acoustic_cost = -decodable->LogLikelihood(frame, arc.ilabel);
183                     double total_cost = tok->cost_ + arc.weight.Value() + acoustic_cost;
184                     if (total_cost >= cutoff) continue;
185                     if (total_cost + beam_ < cutoff) {
186                         cutoff = total_cost + beam_;
187                         Token *new_tok = new Token(arc, acoustic_cost, tok);
188                         unordered_map<StateId, Token*>::iterator find_iter
189                             = cur_toks_.find(arc.nextstate);
190                         if (find_iter == cur_toks_.end()) {
191                             cur_toks_[arc.nextstate] = new_tok;
192                         } else {
193                             if ( *(find_iter->second) < *new_tok ) {
194                                 Token::TokenDelete(find_iter->second);
195                                 find_iter->second = new_tok;
196                             } else {
197                                 Token::TokenDelete(new_tok);
198                             }
199                         }
200                     }
201                 }
202             }
203         }
204     }
205     num_frames_decoded_++;
206 }

```

遍历上一帧所有的stateid-token映射表

获取token对应的state id和token

遍历从每个状态出发的所有弧

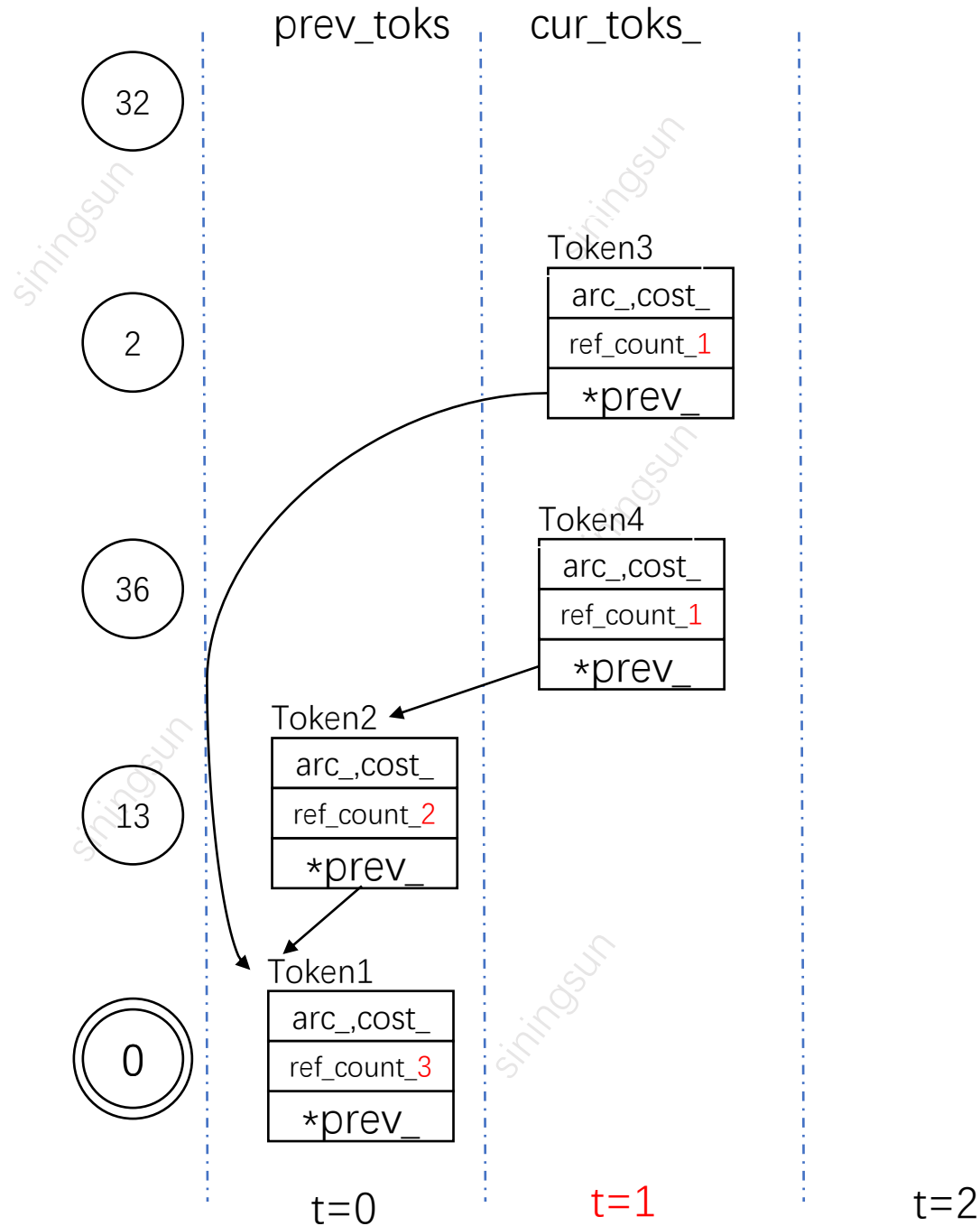
只处理输入为非空的转移, 空转移交给ProcessNonEmitting处理.

原始token中的得分和新转移的声学、图上的得分累计

超过cutoff门限就无需进行传递

如果当前cost+beam比cutoff更小, 那么更新cutoff为一个更加严格的下限

新建一个Token, 确定当前stateid是否存在, 如果不存在(192-194), 直接将新的token加入. 如果存在, 再判断新的token和已经存在的token哪个更好. 如果new_tok比现有的更好, 那么更新当前stateid的Token为new_tok, 并且递归回溯删除原来的token(195-198); 如果当前token不如先有token, 直接删除(199)



逐步分解SimpleDecoder的解码过程

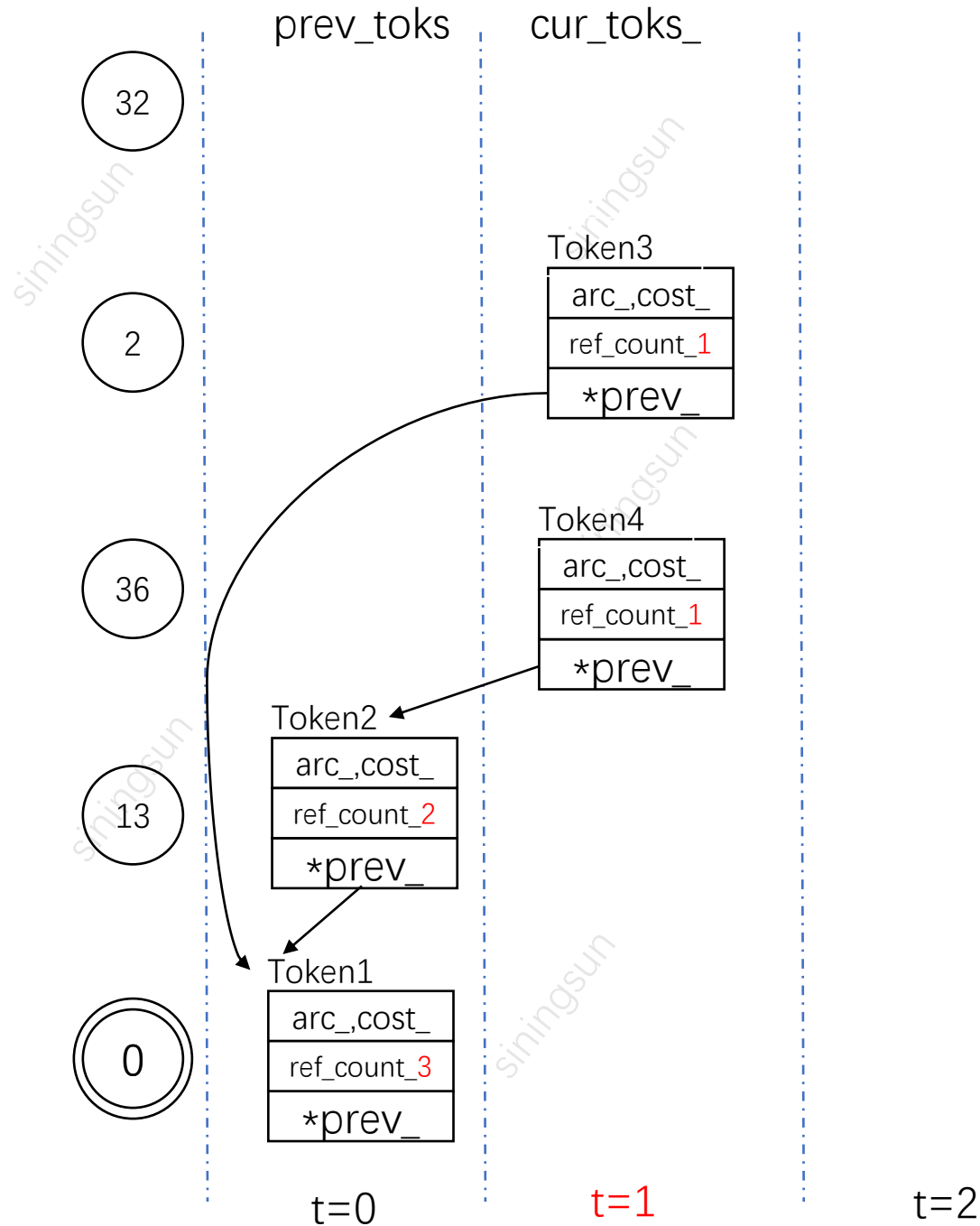
AdvanceDecoding

1. 返回到AdvanceDecoding里面，继续第71行处理非发射弧。
2. 比例中`cur_toks_`里的stateid是2, 36，都没有非发射弧，因此这一步没有进行任务操作
3. 最后第72行进行剪枝，

```

52 void SimpleDecoder::AdvanceDecoding(DecodableInterface *decodable,
53                                     int32 max_num_frames) {
54     KALDI_ASSERT(num_frames_decoded_ >= 0 &&
55                  "You must call InitDecoding() before AdvanceDecoding()");
56     int32 num_frames_ready = decodable->NumFramesReady();
57     // num_frames_ready must be >= num_frames_decoded, or else
58     // the number of frames ready must have decreased (which doesn't
59     // make sense) or the decodable object changed between calls
60     // (which isn't allowed).
61     KALDI_ASSERT(num_frames_ready >= num_frames_decoded_);
62     int32 target_frames_decoded = num_frames_ready;
63     if (max_num_frames >= 0)
64         target_frames_decoded = std::min(target_frames_decoded,
65                                         num_frames_decoded_ + max_num_frames);
66     while (num_frames_decoded_ < target_frames_decoded) {
67         // note: ProcessEmitting() increments num_frames_decoded_
68         ClearToks(prev_toks_);
69         cur_toks_.swap(prev_toks_);
70         ProcessEmitting(decodable);
71         ProcessNonemitting();
72         PruneToks(beam_, &cur_toks_);
73     }
74 }
75

```



逐步分解SimpleDecoder的解码过程

PruneToks()

- SimpleDecoder里面的剪枝操作特别简单，FasterDecoder稍微复杂一下，引入了max和min的active token，并动态调整beam
- 本例子中，我们认为所有的token都在beam以内，没有进行剪枝，`cur_toks_`没有变化，只是对代码进行简单注释

```

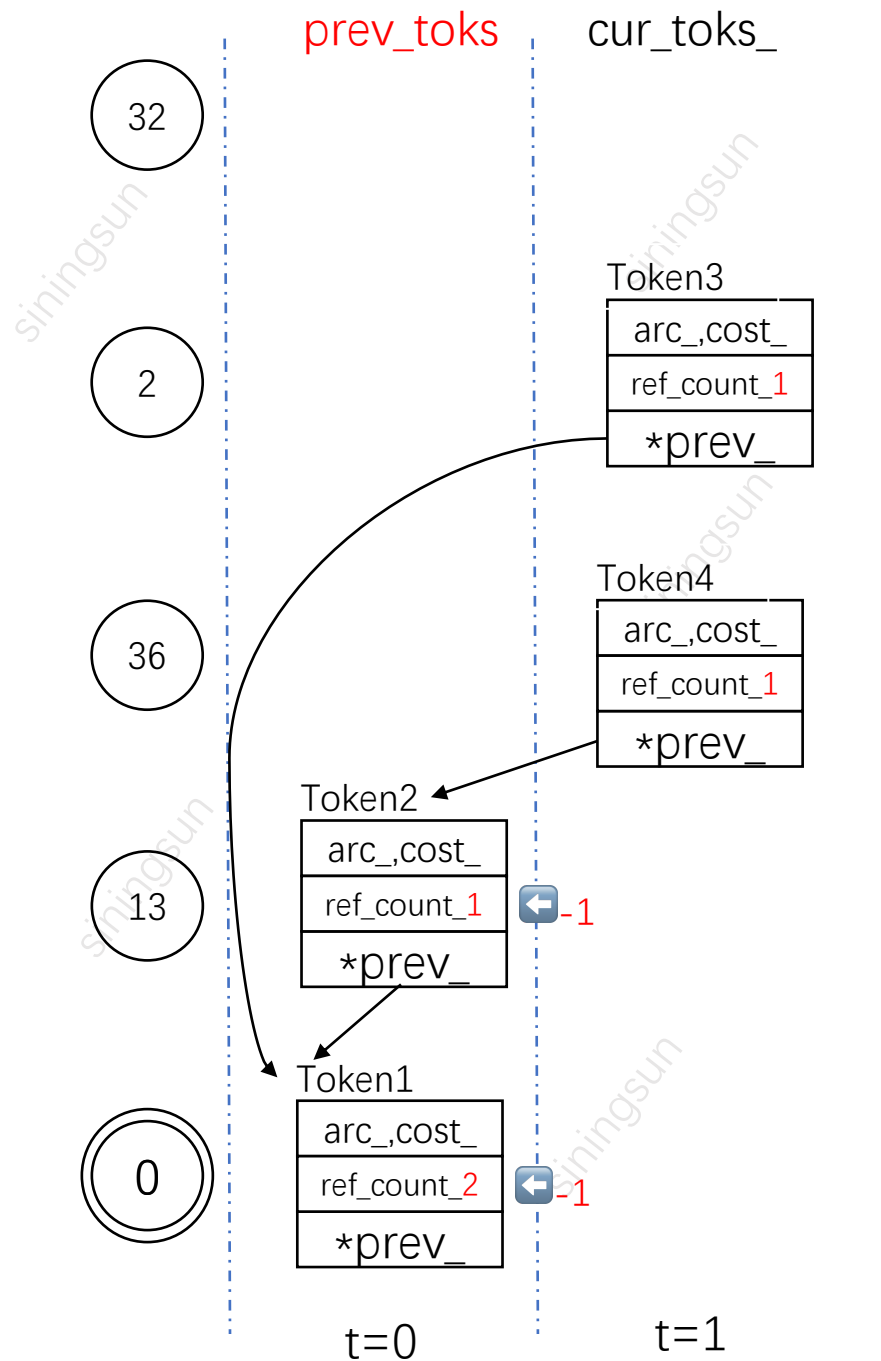
266 // static
267 void SimpleDecoder::PruneToks(BaseFloat beam, unordered_map<StateId, Token*> *toks) {
268     if (toks->empty()) {
269         KALDI_VLOG(2) << "No tokens to prune.\n";
270         return;
271     }
272     double best_cost = std::numeric_limits<double>::infinity();
273     for (unordered_map<StateId, Token*>::iterator iter = toks->begin();
274          iter != toks->end(); ++iter) {
275         best_cost = std::min(best_cost, iter->second->cost_);
276         std::vector<StateId> retained;
277         double cutoff = best_cost + beam;
278         for (unordered_map<StateId, Token*>::iterator iter = toks->begin();
279              iter != toks->end(); ++iter) {
280             if (iter->second->cost_ < cutoff)
281                 retained.push_back(iter->first);
282             else
283                 Token::TokenDelete(iter->second);
284         }
285         unordered_map<StateId, Token*> tmp;
286         for (size_t i = 0; i < retained.size(); i++) {
287             tmp[retained[i]] = (*toks)[retained[i]];
288         }
289         KALDI_VLOG(2) << "Pruned to " << (retained.size()) << " tokens.\n";
290         tmp.swap(*toks);
291     }

```

遍历`cur_toks_`里面所有的元素，或得最优的分数存入`best_cost_`

再一遍遍历`cur_toks_`，如果token得分在beam以内，则将stateid放入`retained`中，否则此token就会被剪枝删除。

将保留的tokens存在tmp中，后面在swap到toks中，也就是`cur_toks_`里面



我们再通过 $t=2$ 的时候跑一遍整个过程，重新回到 AdvanceDecoding 这个函数里面

逐步分解SimpleDecoder的解码过程

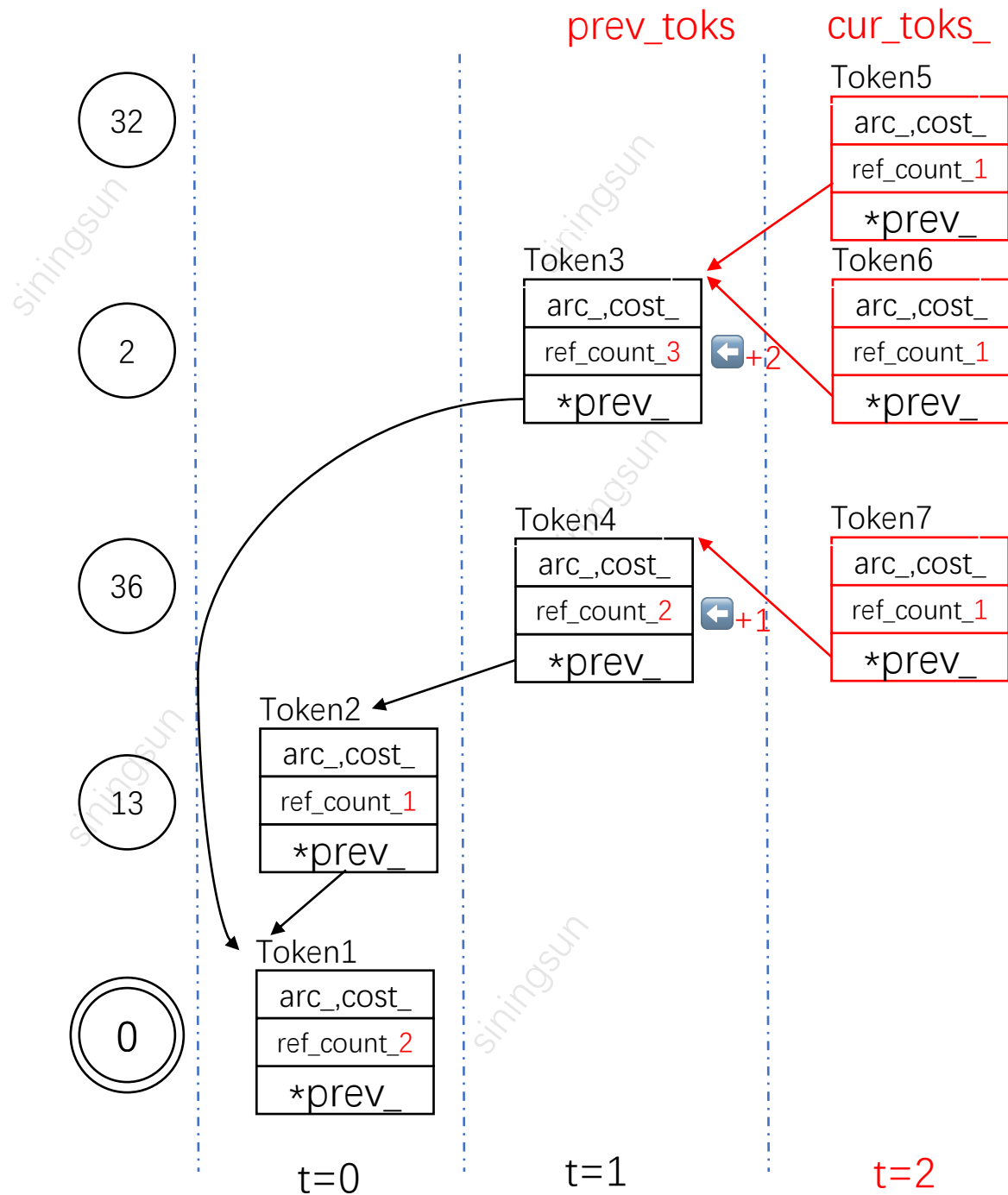
AdvanceDecoding

1. 首先进行ClearToks(prev_toks_)，此时prev_toks_不再是空
2. ClearToks遍历prev_toks_，递归的删除ref_count为1的Token，然后将map清空；注意第261行代码，这里prev_toks_的每个token的ref_count都减去1！这样保证了ref_count就等于回指的Token个数

```

52 void SimpleDecoder::AdvanceDecoding(DecodableInterface *decodable,
53                                     int32 max_num_frames) {
54     KALDI_ASSERT(num_frames_decoded_ >= 0 &&
55                 "You must call InitDecoding() before AdvanceDecoding()");
56     int32 num_frames_ready = decodable->NumFramesReady();
57     // num_frames_ready must be >= num_frames_decoded, or else
58     // the number of frames ready must have decreased (which doesn't
59     // make sense) or the decodable object changed between calls
60     // (which isn't allowed).
61     KALDI_ASSERT(num_frames_ready >= num_frames_decoded_);
62     int32 target_frames_decoded = num_frames_ready;
63     if (max_num_frames >= 0)
64         target_frames_decoded = std::min(target_frames_decoded,
65                                         num_frames_decoded_ + max_num_frames);
66     while (num_frames_decoded_ < target_frames_decoded) {
67         // note: ProcessEmitting() increments num_frames_decoded_
68         ClearToks(prev_toks_);
69         cur_toks_.swap(prev_toks_);
70         ProcessEmitting(decodable);
71         ProcessNonemitting();
72         PruneToks(beam_, &cur_toks_);
73     }
74 }
75
258 void SimpleDecoder::ClearToks(unordered_map<StateId, Token*> &toks) {
259     for (unordered_map<StateId, Token*>::iterator iter = toks.begin();
260          iter != toks.end(); ++iter) {
261         Token::TokenDelete(iter->second);
262     }
263     toks.clear();
264 }

```



逐步分解SimpleDecoder的解码过程

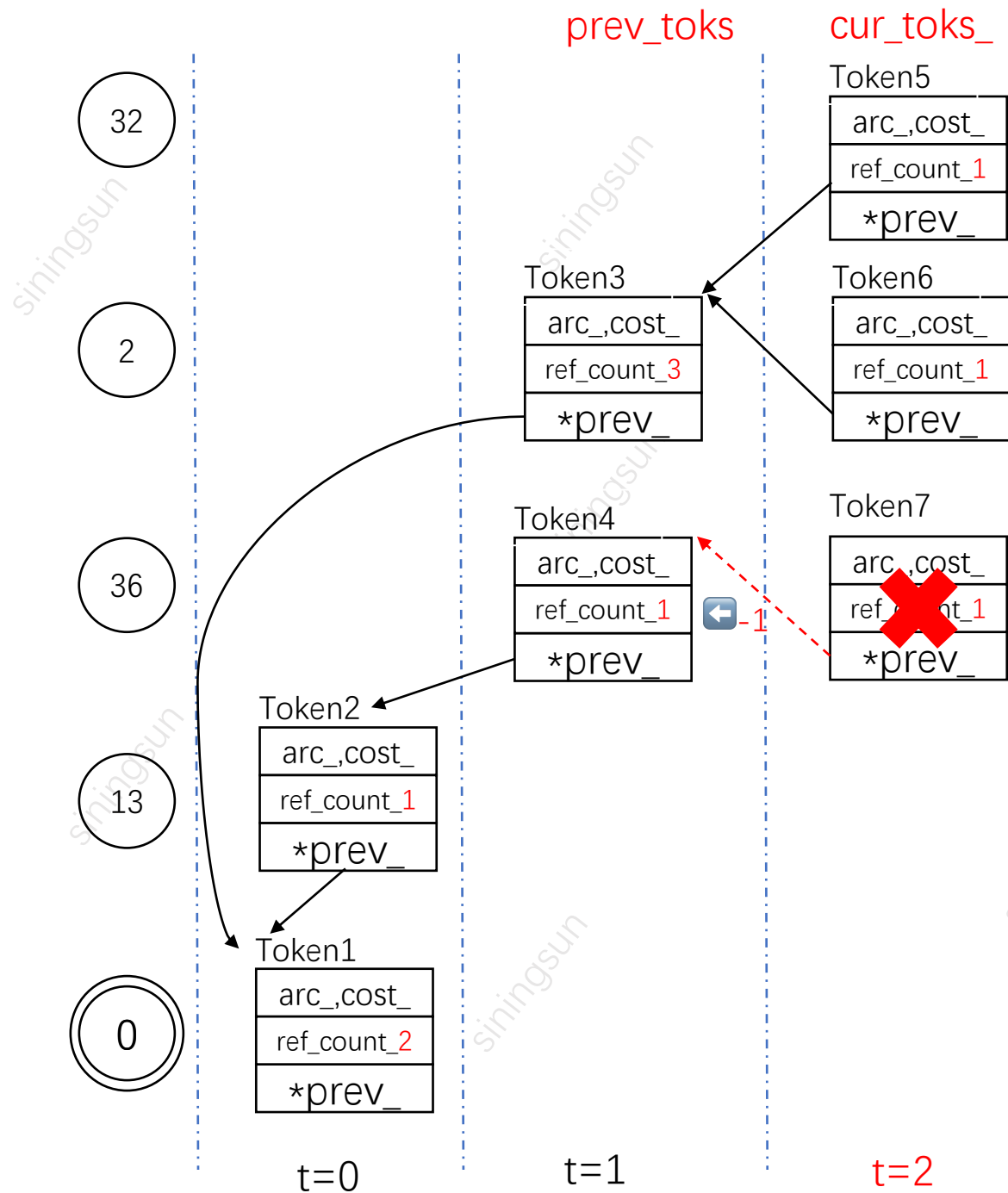
AdvanceDecoding

- 69行, 把`prev_toks_`和`cur_toks_`交换, 此时`prev_toks_`存放的是 $t=1$ 的所有Token, `cur_toks_`被清空, 准备存储 $t=2$ 的Tokens
- 70行, 处理`prev_toks_`里面所有的stateid的发射转移, 此时`prev_toks_`里面只有`<2,Token3>`和`<36,Token4>`,此时可能的转移是 $2 \rightarrow 2$, $2 \rightarrow 32$, $36 \rightarrow 36$, 产生了新的Token5-7
- 我们假设此时没有非发射弧需要处理, 并且假设此时Token7的`cost_`满足剪枝的要求, 然后模拟一下Token7- \rightarrow 4- \rightarrow 2这条路径是如何在 $t=2$ 到 $t=3$ 被剪枝的

```

52 void SimpleDecoder::AdvanceDecoding(DecodableInterface *decodable,
53                                     int32 max_num_frames) {
54     KALDI_ASSERT(num_frames_decoded_ >= 0 &&
55                 "You must call InitDecoding() before AdvanceDecoding()");
56     int32 num_frames_ready = decodable->NumFramesReady();
57     // num_frames_ready must be >= num_frames_decoded, or else
58     // the number of frames ready must have decreased (which doesn't
59     // make sense) or the decodable object changed between calls
60     // (which isn't allowed).
61     KALDI_ASSERT(num_frames_ready >= num_frames_decoded_);
62     int32 target_frames_decoded = num_frames_ready;
63     if (max_num_frames >= 0)
64         target_frames_decoded = std::min(target_frames_decoded,
65                                         num_frames_decoded_ + max_num_frames);
66     while (num_frames_decoded_ < target_frames_decoded) {
67         // note: ProcessEmitting() increments num_frames_decoded_
68         ClearToks(prev_toks_);
69         cur_toks_.swap(prev_toks_);
70         ProcessEmitting(decodable);
71         ProcessNonemitting();
72         PruneToks(beam_, &cur_toks_);
73     }
74 }
75

```



逐步分解SimpleDecoder的解码过程

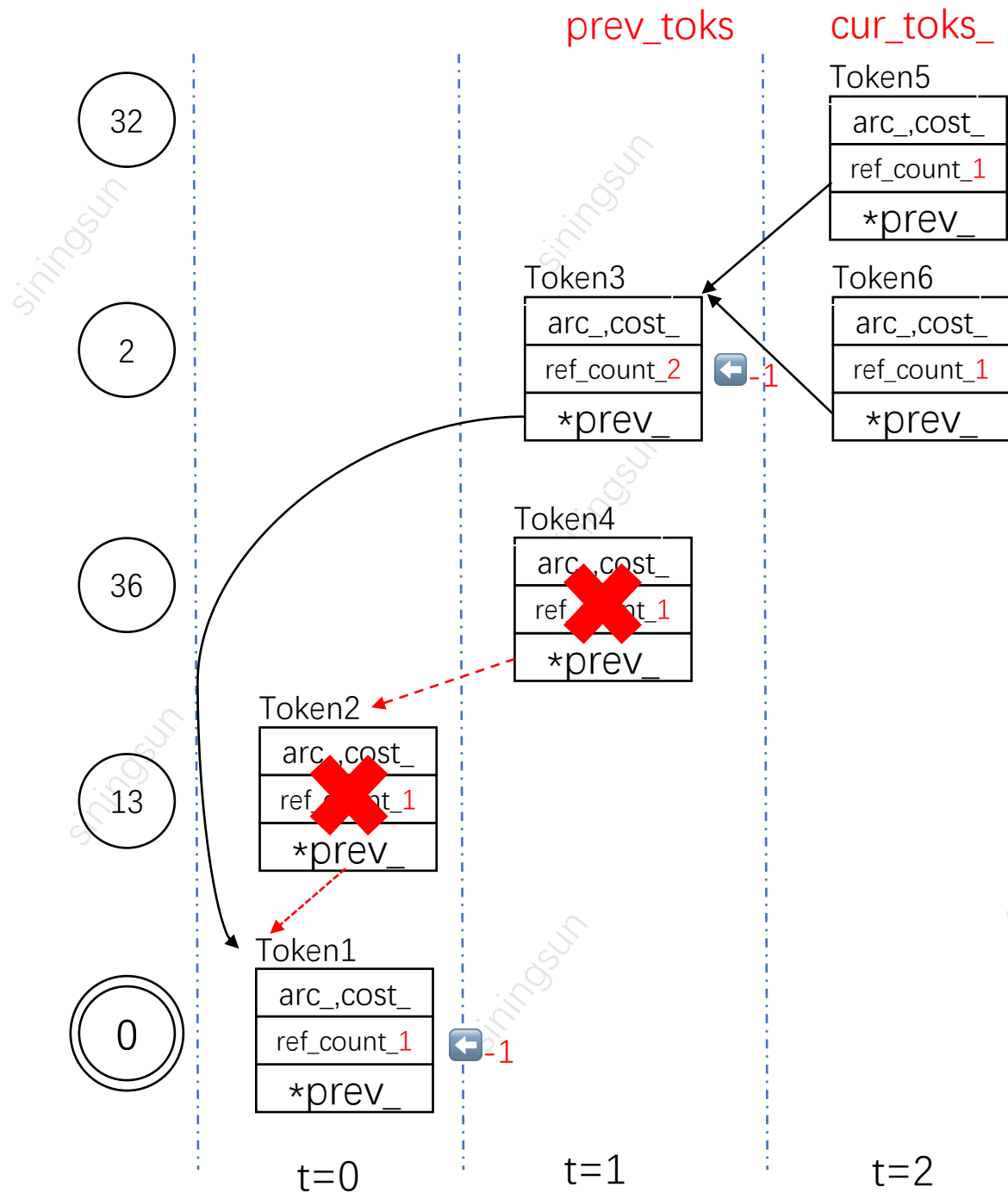
PruneToks()

1. 此时Token7满足282行，调用TokenDelete递归删除
2. 此时Token7的`ref_count`=1,满足while的条件，将Token7的`prev`，也就是Token4赋予tok，删除Token7
3. 重新回到while的条件里之后，tok已经指向了Token4，Token4的`ref_count`-1=1，不满足条件，退出剪枝
4. 此时并没有删除掉Token4，而这一步是在 $t=3$ 完成

```

267 void SimpleDecoder::PruneToks(BaseFloat beam, unordered_map<StateId, Token*> *toks) {
268     if (toks->empty()) {
269         KALDI_VLOG(2) << "No tokens to prune.\n";
270         return;
271     }
272     double best_cost = std::numeric_limits<double>::infinity();
273     for (unordered_map<StateId, Token*>::iterator iter = toks->begin();
274          iter != toks->end(); ++iter)
275         best_cost = std::min(best_cost, iter->second->cost_);
276     std::vector<StateId> retained;
277     double cutoff = best_cost + beam;
278     for (unordered_map<StateId, Token*>::iterator iter = toks->begin();
279          iter != toks->end(); ++iter) {
280         if (iter->second->cost_ < cutoff)
281             retained.push_back(iter->first);
282     }
283     for (StateId state_id : retained) {
284         Token *tok = toks->at(state_id);
285         static void TokenDelete(Token *tok) {
286             while (--tok->ref_count_ == 0) {
287                 Token *prev = tok->prev_;
288                 delete tok;
289                 if (prev == NULL) return;
290                 else tok = prev;
291             }
292         }
293         #ifdef KALDI_PARANOID
294             KALDI_ASSERT(tok->ref_count_ > 0);
295         #endif
296     }
297 }

```

逐步分解SimpleDecoder的解码过程

AdvanceDecoding

- t=3时, 68行首先进行ClearToks(prev_toks_) 此时prev_toks_ 还停留在t=1, 此时对 prev_toks_ 进行clear, 对Token3和4进行递归删除, 此时Token4, Token2都会被删除, Token1和Token3的ref_count_ 会减1.
- 69行将t=2的toks_ 交换给prev_toks_, 继续解码

```

52 void SimpleDecoder::AdvanceDecoding(DecodableInterface *decodable,
53                                     int32 max_num_frames) {
54     KALDI_ASSERT(num_frames_decoded_ >= 0 &&
55                  "You must call InitDecoding() before AdvanceDecoding()");
56     int32 num_frames_ready = decodable->NumFramesReady();
57     // num_frames_ready must be >= num_frames_decoded, or else
58     // the number of frames ready must have decreased (which doesn't
59     // make sense) or the decodable object changed between calls
60     // (which isn't allowed).
61     KALDI_ASSERT(num_frames_ready >= num_frames_decoded_);
62     int32 target_frames_decoded = num_frames_ready;
63     if (max_num_frames >= 0)
64         target_frames_decoded = std::min(target_frames_decoded,
65                                           num_frames_decoded_ + max_num_frames);
66     while (num_frames_decoded_ < target_frames_decoded) {
67         // note: ProcessEmitting() increments num_frames_decoded_
68         ClearToks(prev_toks_);
69         cur_toks_.swap(prev_toks_);
70         ProcessEmitting(decodable);
71         ProcessNonemitting();
72         PruneToks(beam_, &cur_toks_);
73     }
74 }
75
258 void SimpleDecoder::ClearToks(unordered_map<StateId, Token*> &toks) {
259     for (unordered_map<StateId, Token*>::iterator iter = toks.begin();
260          iter != toks.end(); ++iter) {
261         Token::TokenDelete(iter->second);
262     }
263     toks.clear();
264 }

```