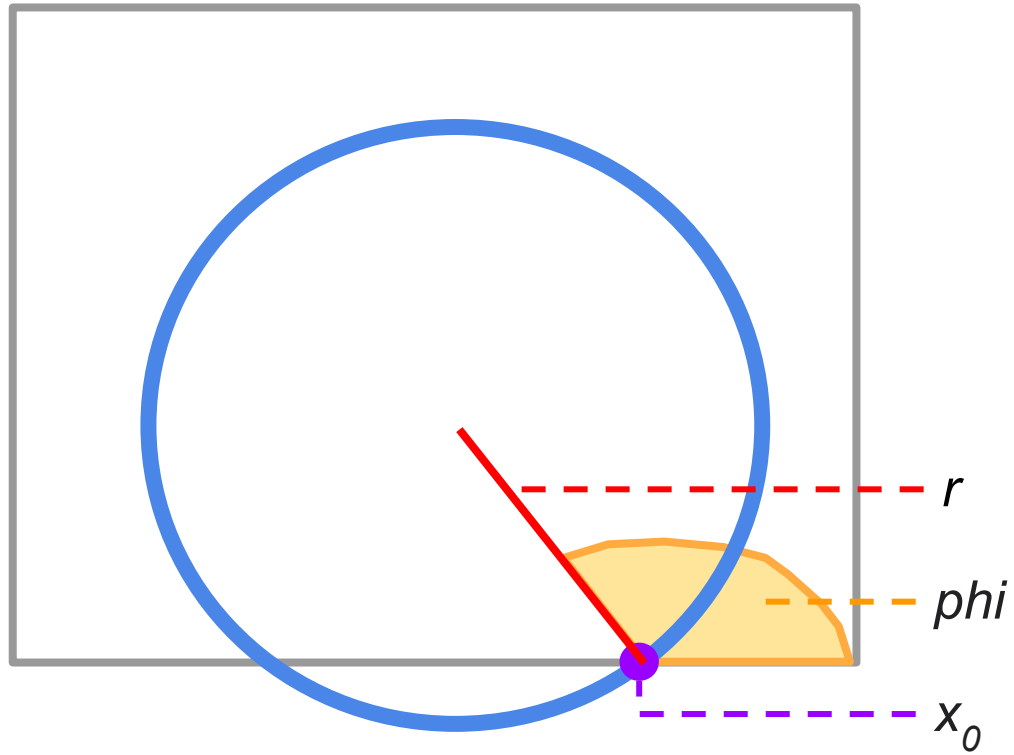# Language Model Training with STT Toy Data Generator
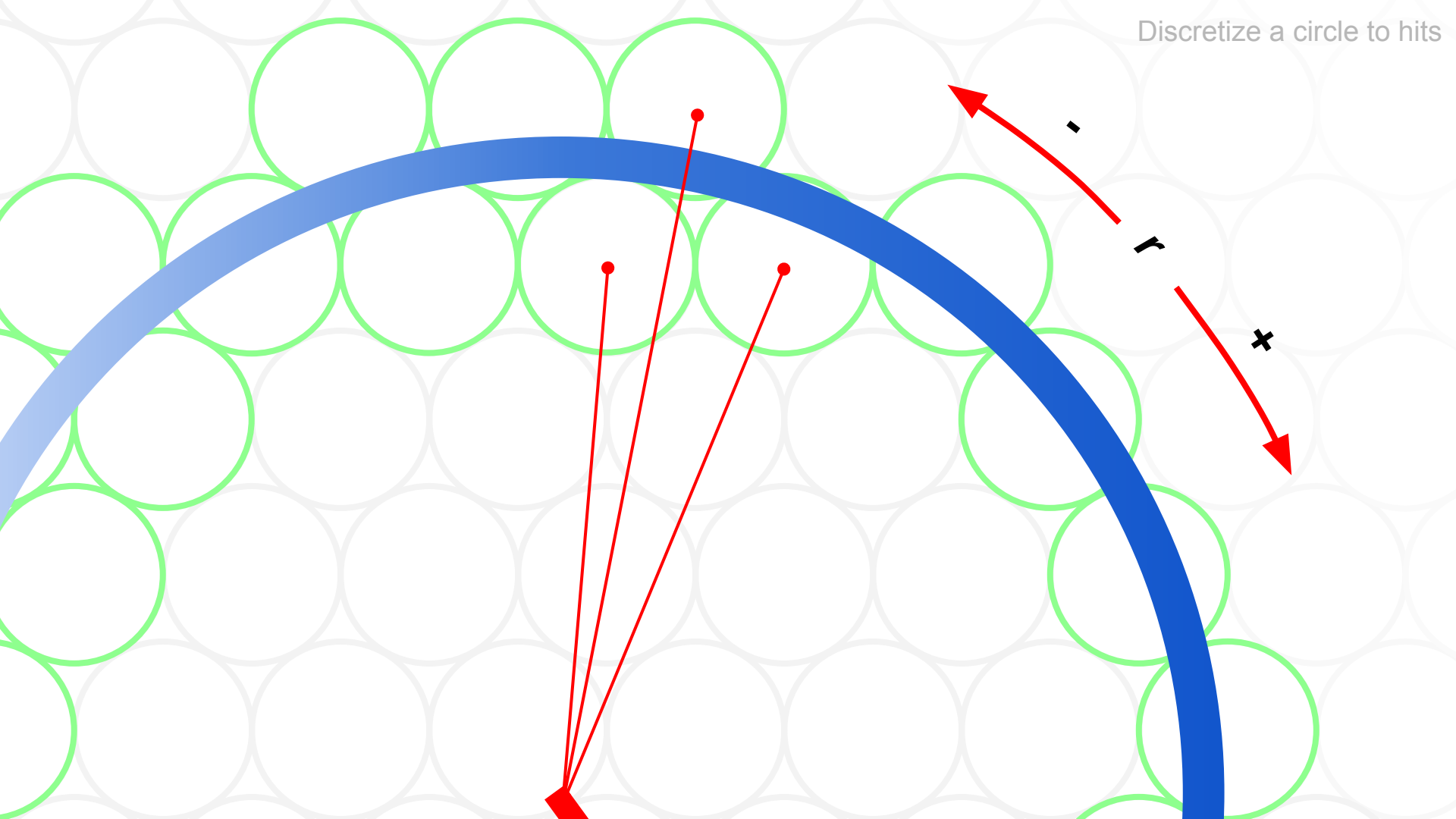
Jakapat Kannika
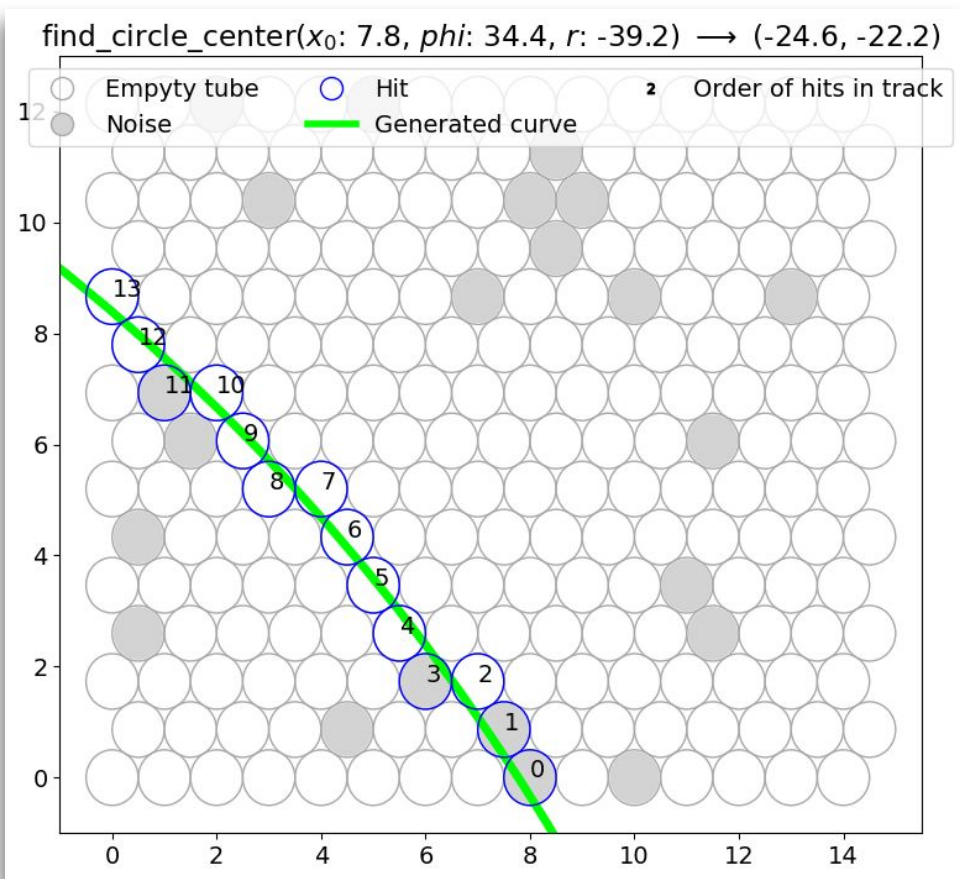
# Toy Data Generator for STT

Simulation frame

Discretize a circle to hits

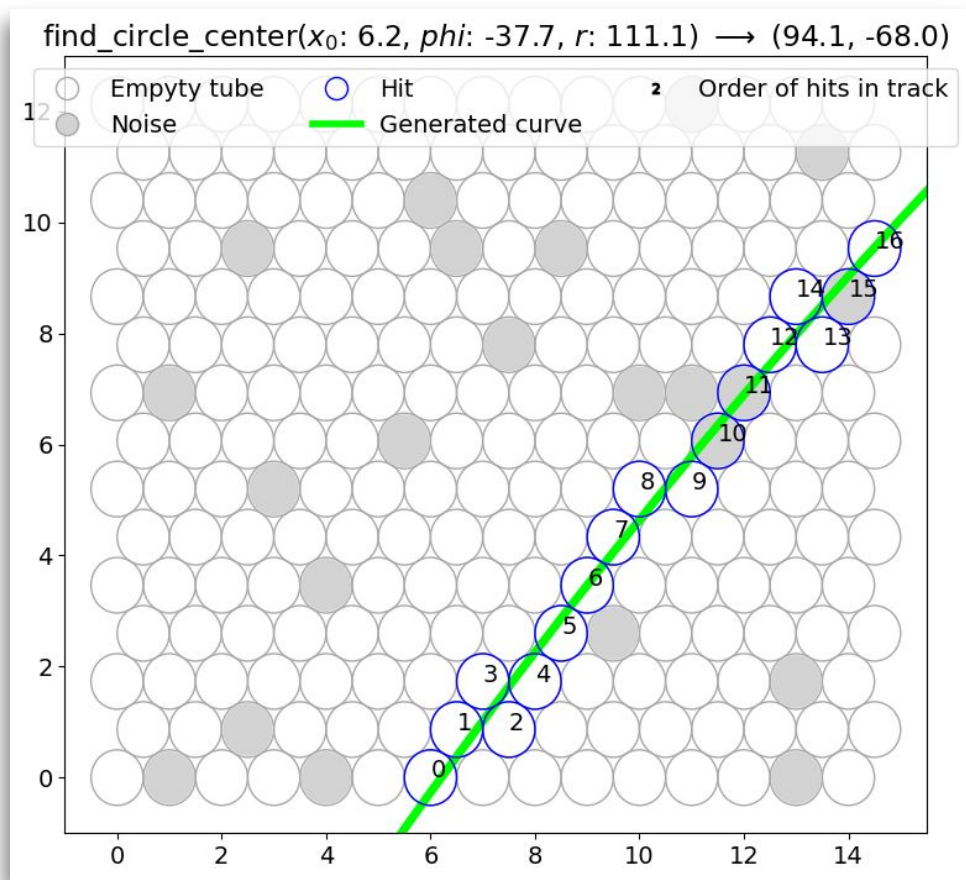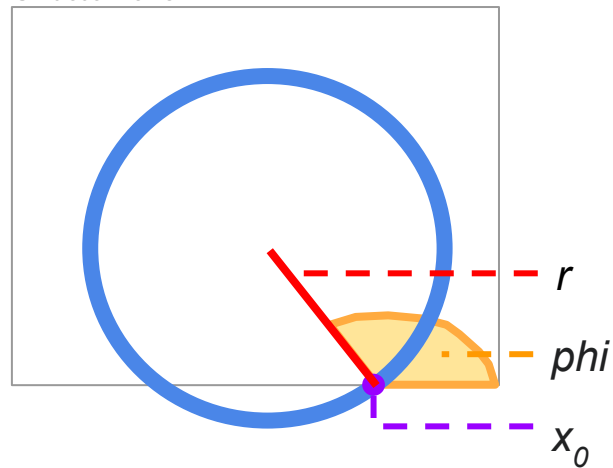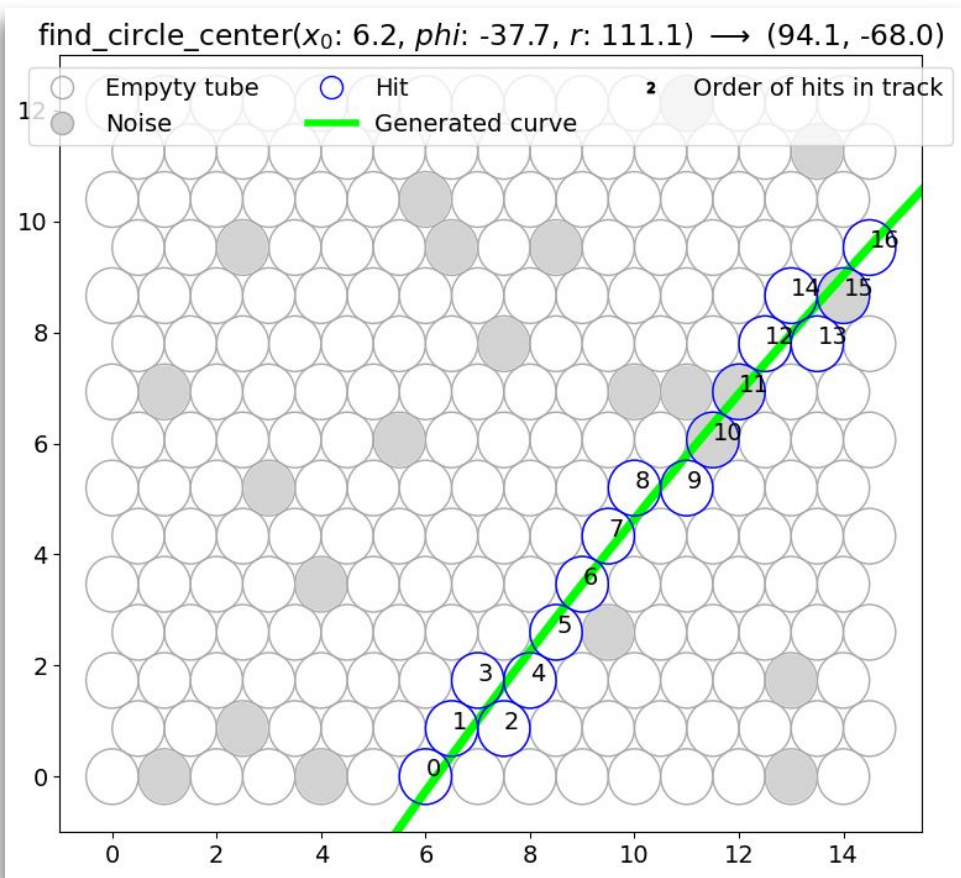find_circle_center($x_0$: 6.2, *phi*: -37.7, *r*: 111.1) ⟶ (94.1, -68.0)

Empty tube   ◯ Hit        **2** Order of hits in track
Noise        ─── Generated curve

Simulation frame

find_circle_center($x_0$: 6.2, *phi*: -37.7, *r*: 111.1) $\longrightarrow$ (94.1, -68.0)

Positions: [
[6.0, 0.0], [6.5, 0.9], [7.5, 0.9], [7.0, 1.7],
[8.0, 1.7], [8.5, 2.6], [9.0, 3.5], [9.5, 4.3],
[10.0, 5.2], [11.0, 5.2], [11.5, 6.1], [12.0,
6.9], [12.5, 7.8], [13.5, 7.8], [13.0, 8.7],
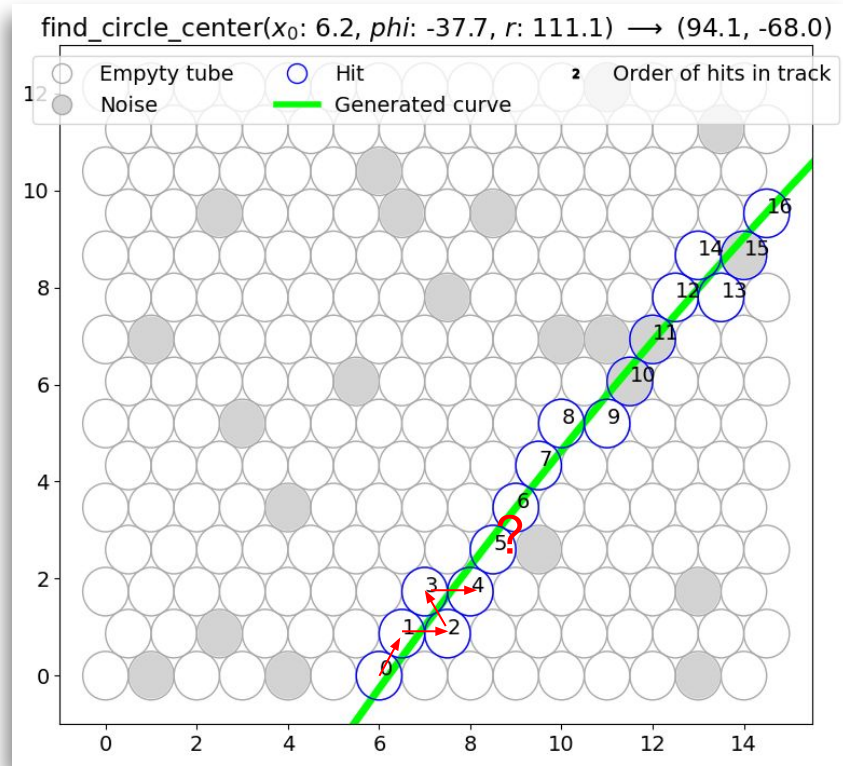[14.0, 8.7], [14.5, 9.5]]

**Tracking features**

Moving directions: [
60, 0, 120, 0, 60, 60, 60, 60, 0, 60, 60,
60, 0, 120, 0, 60]

Neighbor patterns: [
 [1, 41, 7, 56, 13, 41, 25, 9, 40, 5, 11,
13, 41, 7, 56, 13, 8]

# Language Model

find_circle_center($x_0$: 6.2, $phi$: -37.7, $r$: 111.1) ⟶ (94.1, -68.0)

A statistical language model is a probability distribution over sequences of words. Given such a sequence, say of length m, it assigns a probability $P(w_1, \ldots, w_m)$ to the whole sequence.*
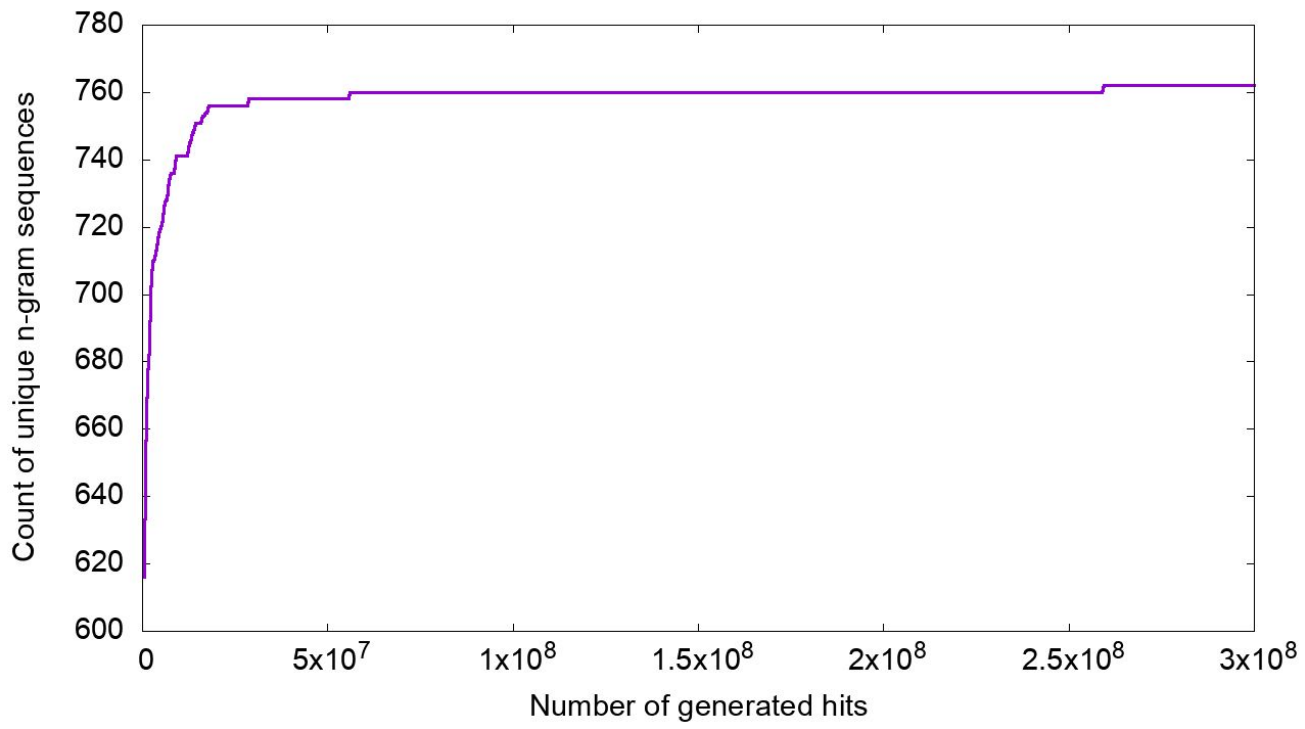
…

# Moving directions:

60, 0, 120, 0, 60, 60, 60, 60, 0, 60, 60, 60, 0, 120, 0, 60

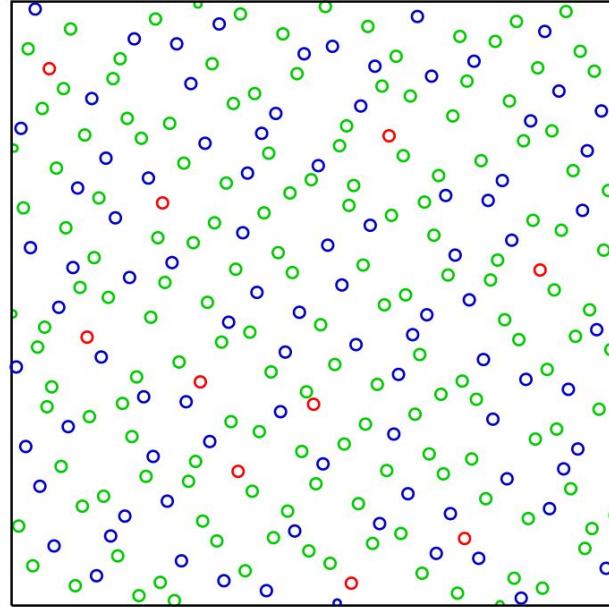| Pattern | Count | Prob. |
|---|---|---|
| 60, 0, 120, 0 | 2 | 1.00 |
| 0, 120, 0, 60 | 2 | 1.00 |
| 120, 0, 60, 60 | 1 | 1.00 |
| 0, 60, 60, 60 | 2 | 1.00 |
| 60, 60, 60, 60 | 1 | 0.33 |
| 60, 60, 60, 0 | 2 | 0.66 |
| 60, 60, 0, 60 | 1 | 0.50 |
| 60, 0, 60, 60 | 1 | 1.00 |
| 60, 60, 0, 120 | 1 | 0.50 |

# Current training models

- Training feature: moving directions,
- Language models: 5-gram, 10-gram, 15-gram models,
- Sizes of simulation frames: 15 x 15, 20 x 20, 25 x 25 tubes
- Noise: 0 noise hit.
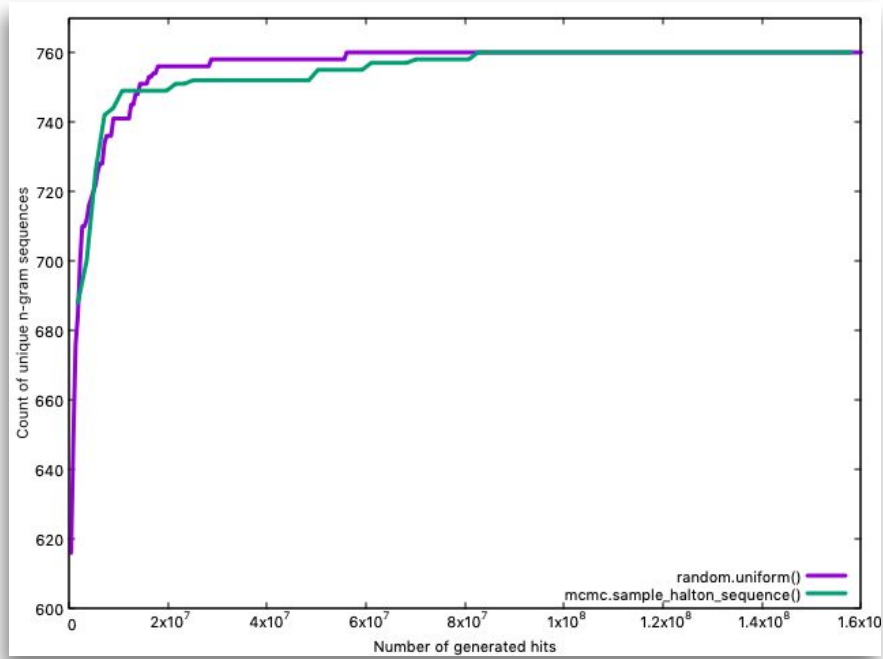
# Optimize training speed using halton sequence
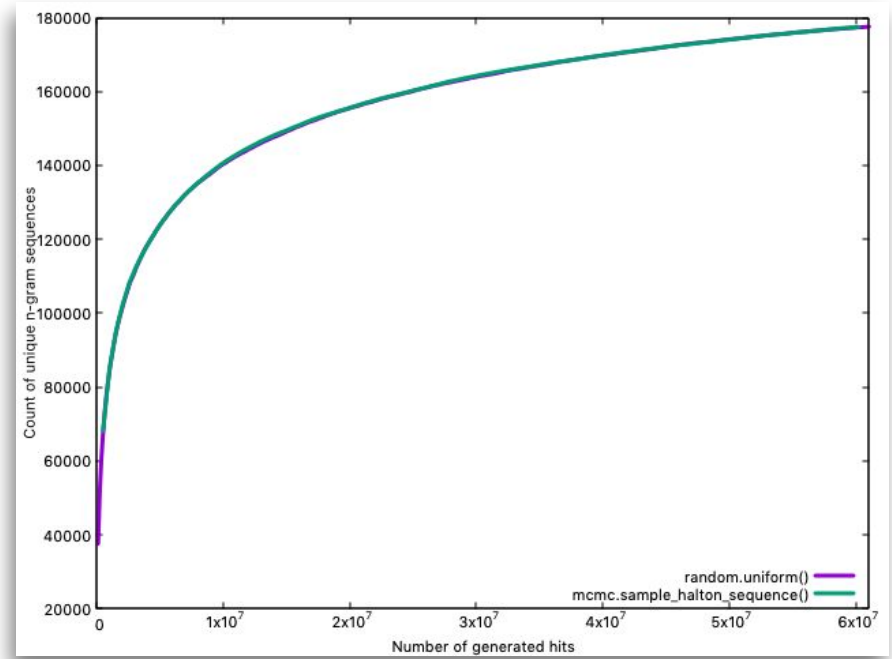


Pseudorandom

Halton sequence

5-gram model



10-gram model

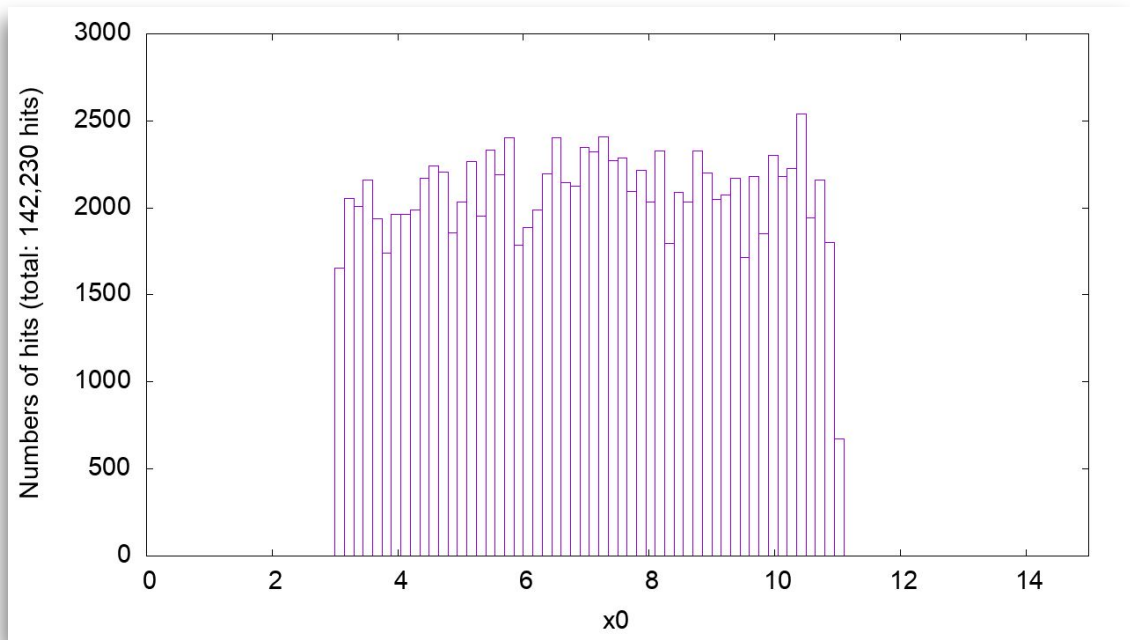# Check distributions of hits

Simulation frame:

- Width = 15 tubes,
- Height = 15 rows.

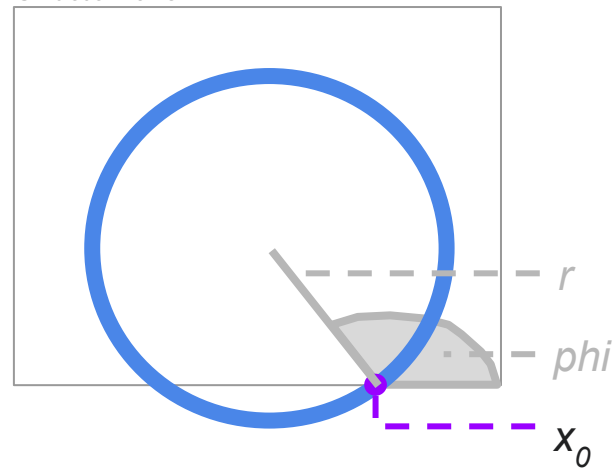Training language model:

- 5-gram model for moving directions.
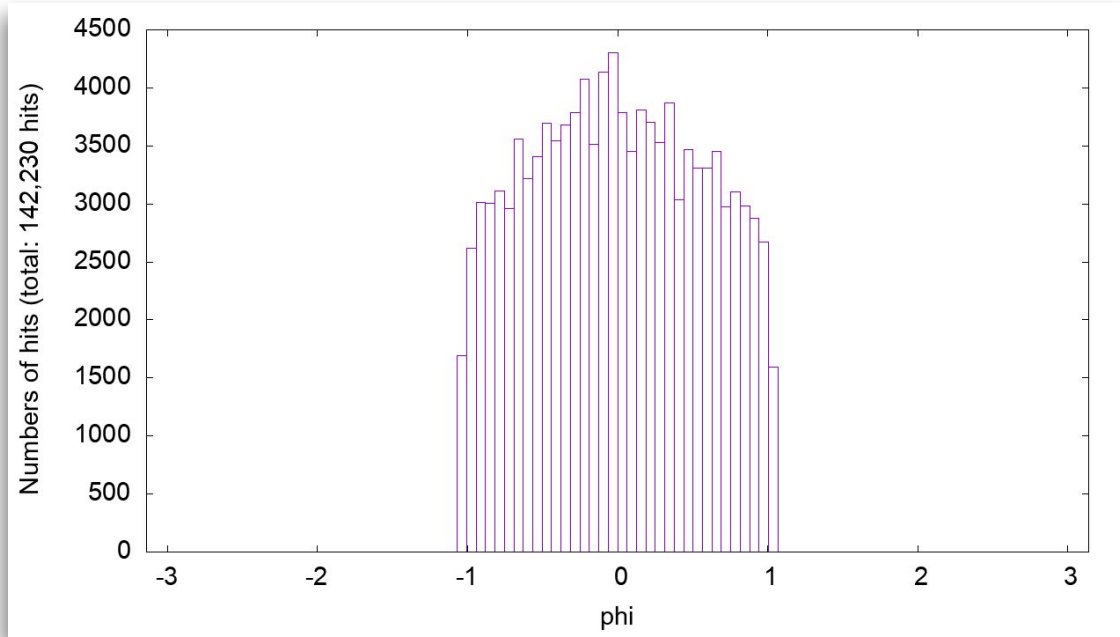
Number of generating data:
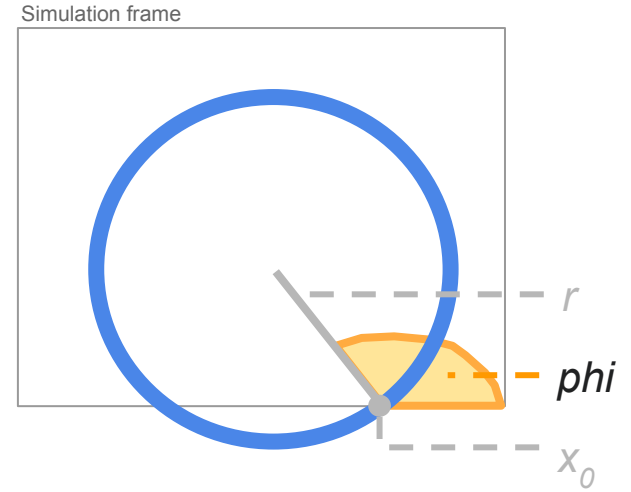
- 142,260 hits (10,000 tracks)
- 0 noise hit.

x0 = random.uniform(3, 11)
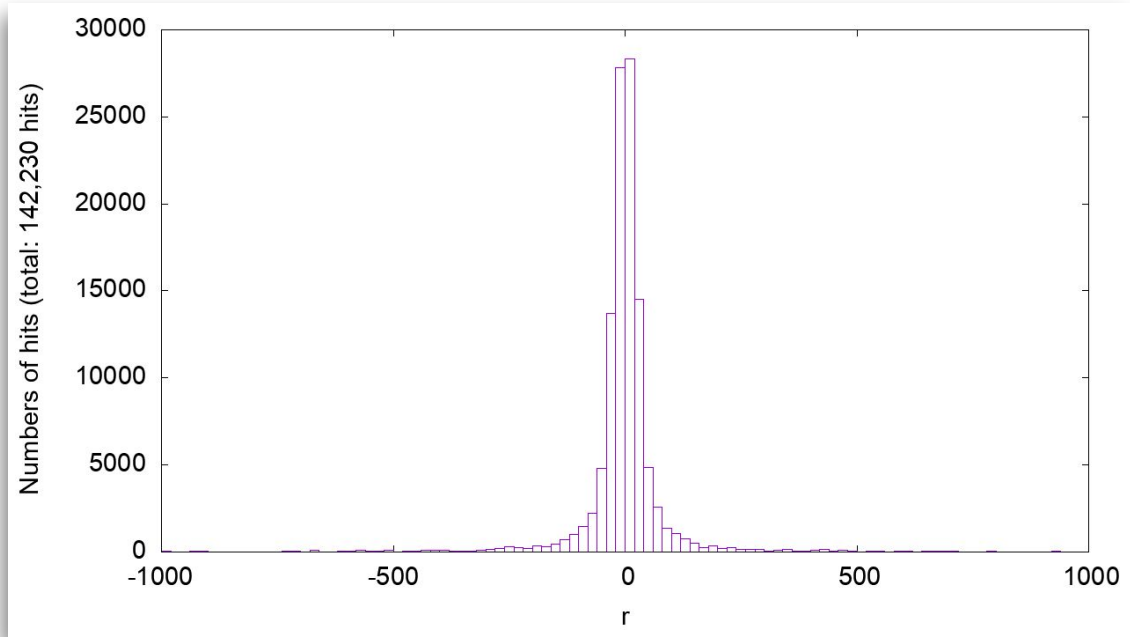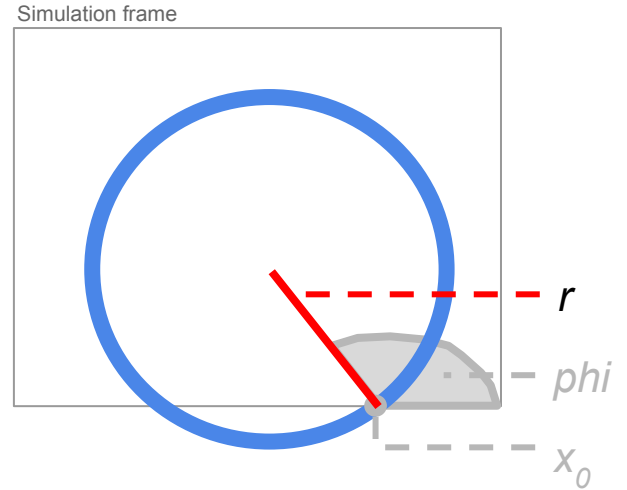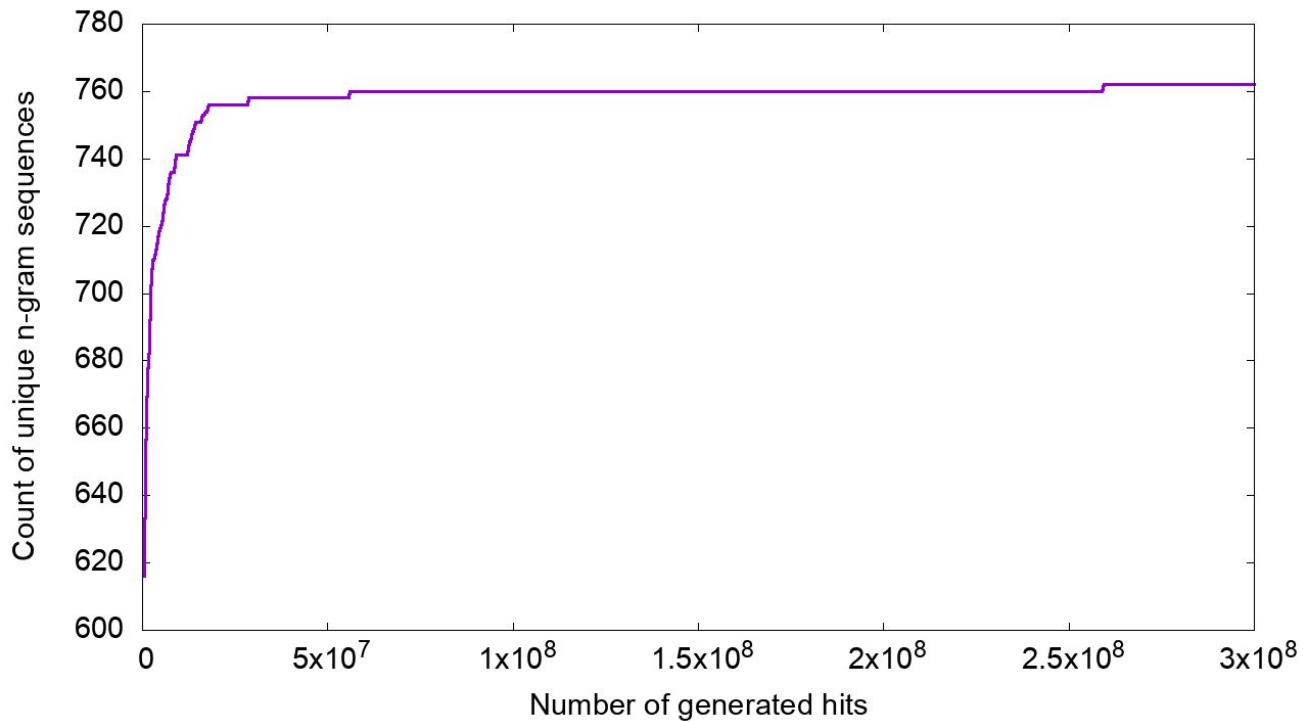
phi = random.uniform(-1 * math.pi / 3.0, math.pi / 3.0)

a = random.uniform(0.001, 0.1)
r = random.choice([-1, 1]) * (1 / a)

# Check distribution of new patterns
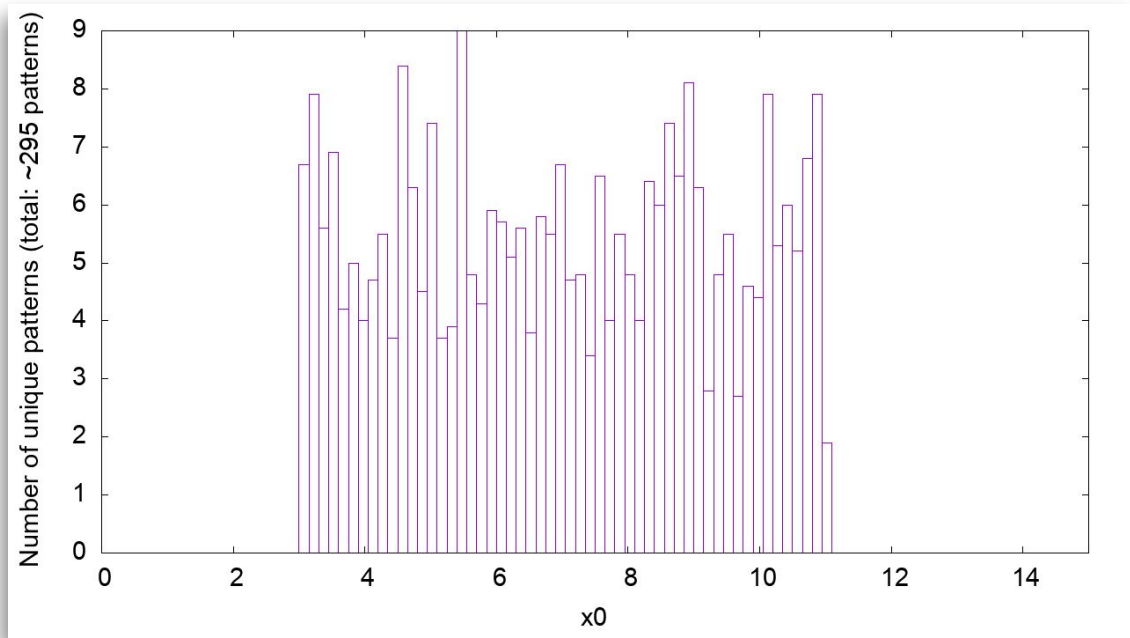
Simulation frame:

- Width = 15 tubes,
- Height = 15 rows.

Training language model:

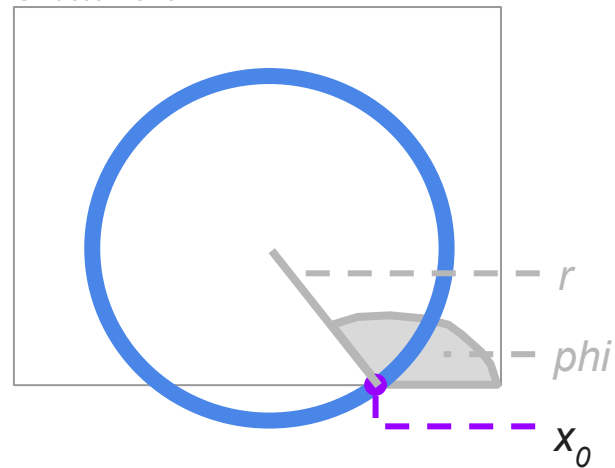- 5-gram model for moving directions.

Number of generating data:

- ~141,982 hits (10,000 tracks)
- 0 noise hit.

x0 = random.uniform(3, 11)

phi = random.uniform(-1 * math.pi / 3.0, math.pi / 3.0)

r = random.choice([-1, 1]) * random.uniform(10, 1000)

a = random.uniform(0.001, 0.1)
r = random.choice([-1, 1]) * (1 / a)



r = random.choice([-1, 1]) * random.uniform(10, 1000)

# Check for a bottleneck in the data generation

Simulation frame:

- Width = 15 tubes,
- Height = 15 rows.

Training language model:
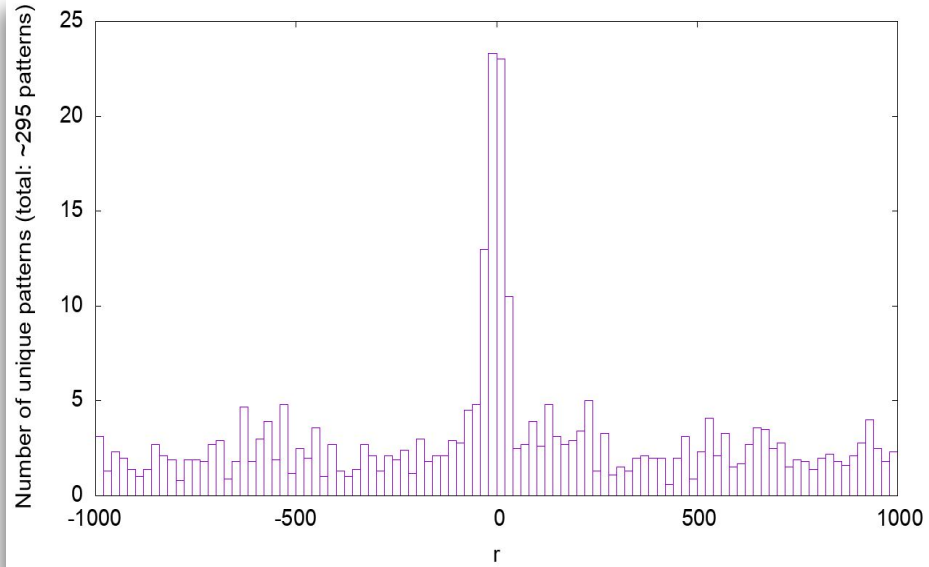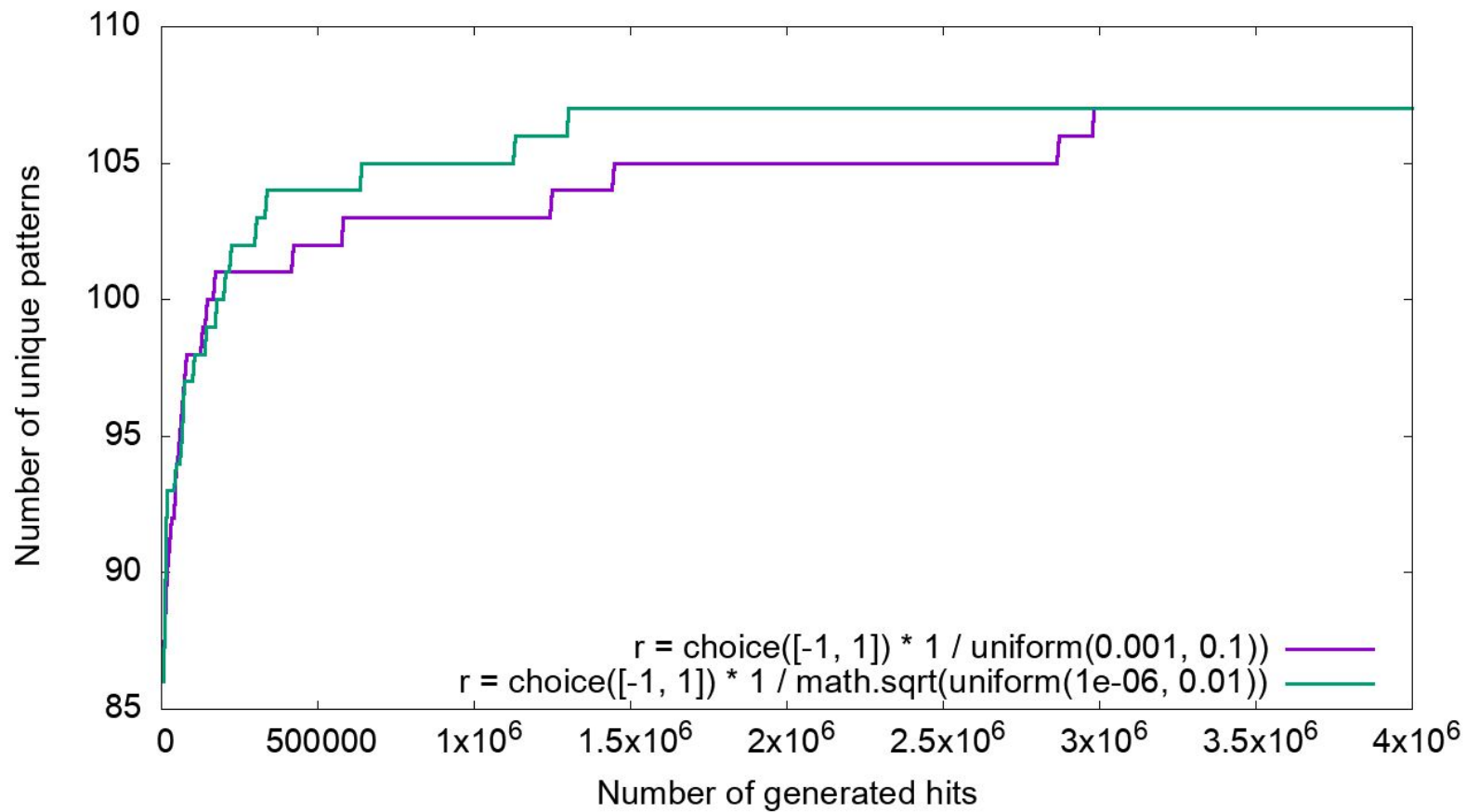
- 3-gram model for moving directions.

Number of generating data:

- ~4,000,000 hits
- 0 noise hit.

# Summary and outlook

Summary:

- The new toy data generator can generate data with the geometry similar to the STT,
- The generator can produce consistent patterns that can be used in language model training,
- Feature extractors for moving directions and neighbor patterns are available for the new geometry,
- Slow in speed of training could be caused by a bottleneck in the data generation.

Outlook:

- Finish language model training for moving directions and neighbor patterns,
- Implement isochrone radius for the new data generator.