

Montreal Transit Summary

March 9, 2026

1 MTL Transit Transformer Analysis

1.1 Overview

This project trains a GPT-style transformer on **STM (Societe de transport de Montreal) GTFS data** to learn the structural grammar of bus and metro trips.

Each trip is converted into a symbolic token sequence:

```
<BOS> <ROUTE_80> <DIR_1> <SERVICE_weekday> <STOP_5241> <SEQ_01> <T_482> ... <EOS>
```

The trained model learns route topology, stop ordering, and schedule rhythm — enabling structural analysis of the transit network through language modeling techniques.

Key analyses: - Next-token prediction accuracy - Corruption detection (topology validation) - Route embedding clustering - Network bottleneck identification - Schedule padding detection - Route variant discovery

```
[1]: import os
import json
import pickle
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import matplotlib
matplotlib.rcParams['figure.dpi'] = 120
matplotlib.rcParams['figure.figsize'] = (10, 5)

ROOT = os.path.abspath(os.path.join(os.getcwd(), '..'))
DATA_DIR = os.path.join(ROOT, 'data', 'mtl_transit')
ANALYSIS_DIR = os.path.join(ROOT, 'analysis')
```

1.2 1. Dataset Overview

```
[2]: with open(os.path.join(DATA_DIR, 'meta.pkl'), 'rb') as f:
    meta = pickle.load(f)

with open(os.path.join(DATA_DIR, 'input.txt'), 'r') as f:
    lines = f.readlines()
```

```

token_counts = [len(line.strip().split()) for line in lines]
train_ids = np.memmap(os.path.join(DATA_DIR, 'train.bin'), dtype=np.uint16,
    ↪mode='r')
val_ids = np.memmap(os.path.join(DATA_DIR, 'val.bin'), dtype=np.uint16, mode='r')

print(f"Trips:           {len(lines):,}")
print(f"Vocabulary size:  {meta['vocab_size']:,}")
print(f"Avg trip length:    {np.mean(token_counts):.1f} tokens")
print(f"Max trip length:    {max(token_counts)} tokens")
print(f"Train tokens:       {len(train_ids):,}")
print(f"Val tokens:         {len(val_ids):,}")
print(f"Total tokens:       {len(train_ids) + len(val_ids):,}")

```

```

Trips:           273,967
Vocabulary size: 10,647
Avg trip length: 114.6 tokens
Max trip length: 356 tokens
Train tokens:    28,695,894
Val tokens:      2,706,326
Total tokens:    31,402,220

```

```

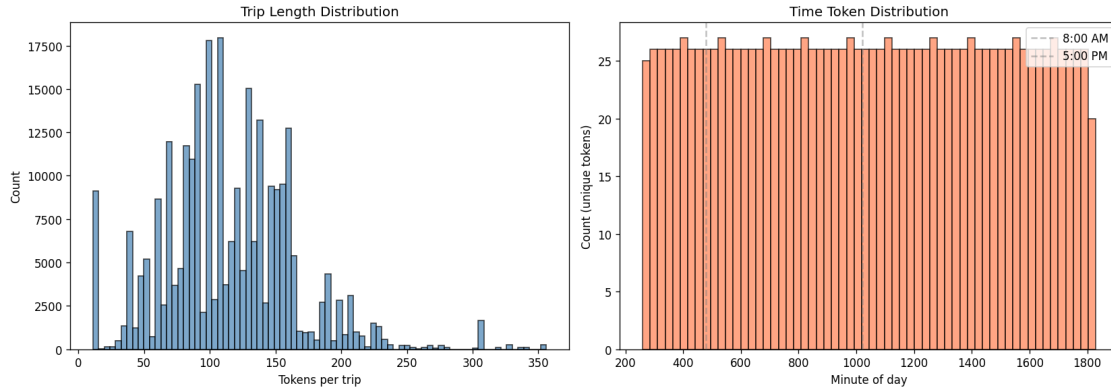
[3]: fig, axes = plt.subplots(1, 2, figsize=(14, 5))

axes[0].hist(token_counts, bins=80, edgecolor='black', alpha=0.7,
    ↪color='steelblue')
axes[0].set_xlabel('Tokens per trip')
axes[0].set_ylabel('Count')
axes[0].set_title('Trip Length Distribution')

time_tokens = [t for t in meta['stoi'] if t.startswith('<T_')]
time_values = sorted([int(t.split('_')[1].rstrip('>')) for t in time_tokens])
axes[1].hist(time_values, bins=60, edgecolor='black', alpha=0.7, color='coral')
axes[1].set_xlabel('Minute of day')
axes[1].set_ylabel('Count (unique tokens)')
axes[1].set_title('Time Token Distribution')
axes[1].axvline(x=480, color='gray', linestyle='--', alpha=0.5, label='8:00 AM')
axes[1].axvline(x=1020, color='gray', linestyle='--', alpha=0.5, label='5:00 PM')
axes[1].legend()

plt.tight_layout()
plt.show()

```



1.3 2. Example Trip Sequences

Each trip is encoded as a single line of symbolic tokens. The structure follows:

Token	Meaning
<BOS> / <EOS>	Trip boundaries
<ROUTE_x>	Route ID
<DIR_x>	Direction (0 or 1)
<SERVICE_x>	Day type: weekday, saturday, or sunday
<STOP_x>	Stop ID
<SEQ_xx>	Stop sequence number
<T_x>	Departure minute of day

```
[4]: for i in [0, 100, 5000]:
    tokens = lines[i].strip().split()
    route = next(t for t in tokens if t.startswith('<ROUTE_'))
    n_stops = sum(1 for t in tokens if t.startswith('<STOP_'))
    print(f"Trip {i}: {route}, {n_stops} stops")
    print(f" {' '.join(tokens[:16])} ...")
    print()
```

Trip 0: <ROUTE_11>, 25 stops

```
<BOS> <ROUTE_11> <DIR_1> <SERVICE_weekday> <STOP_52356> <SEQ_01> <T_456>
<STOP_52271> <SEQ_02> <T_456> <STOP_52160> <SEQ_03> <T_457> <STOP_52137>
<SEQ_04> <T_458> ...
```

Trip 100: <ROUTE_51>, 51 stops

```
<BOS> <ROUTE_51> <DIR_0> <SERVICE_weekday> <STOP_50110> <SEQ_01> <T_611>
<STOP_53792> <SEQ_02> <T_612> <STOP_50766> <SEQ_03> <T_613> <STOP_50716>
<SEQ_04> <T_614> ...
```

Trip 5000: <ROUTE_491>, 8 stops

```
<BOS> <ROUTE_491> <DIR_0> <SERVICE_weekday> <STOP_60413> <SEQ_01> <T_727>
```

```
<STOP_57198> <SEQ_02> <T_728> <STOP_62120> <SEQ_03> <T_729> <STOP_57214>  
<SEQ_04> <T_731> ...
```

1.4 3. Model Performance

A 6-layer transformer (6 heads, 384 embedding dim, ~15M parameters) was trained on 28.7M tokens for 6000 iterations.

1.4.1 Next-Token Prediction

The model was evaluated on validation trips at multiple prefix lengths. Results demonstrate strong learning of trip grammar.

```
[5]: eval_path = os.path.join(ANALYSIS_DIR, 'next_token_eval.json')  
with open(eval_path, 'r') as f:  
    eval_results = json.load(f)  
  
if 'overall' in eval_results:  
    print(f"Overall top-1 accuracy: {eval_results['overall']['top1']:.1%}")  
    print(f"Overall top-5 accuracy: {eval_results['overall']['top5']:.1%}")  
    print()  
  
if 'next_stop' in eval_results:  
    ns = eval_results['next_stop']  
    print(f"Next stop top-1: {ns['top1']:.1%}")  
    print(f"Next stop top-5: {ns['top5']:.1%}")  
    print()  
  
if 'next_time' in eval_results:  
    nt = eval_results['next_time']  
    print(f"Next time exact match: {nt['exact_match']:.1%}")  
    print(f"Next time MAE: {nt['mean_abs_error_min']:.1f} minutes")  
    print()  
  
if 'by_class' in eval_results:  
    print("By token class:")  
    print(f"  {'Class':>10s}  {'Top-1':>8s}  {'Top-5':>8s}")  
    for cls in ['ROUTE', 'DIR', 'SERVICE', 'STOP', 'SEQ', 'TIME']:  
        if cls in eval_results['by_class']:  
            d = eval_results['by_class'][cls]  
            print(f"  {cls:>10s}  {d['top1']:>8.1%}  {d['top5']:>8.1%}")
```

Overall top-1 accuracy: 99.4%

Overall top-5 accuracy: 99.4%

Next stop top-1: 98.8%

Next stop top-5: 99.1%

Next time exact match: 86.3%

Next time MAE: 13.4 minutes

By token class:

Class	Top-1	Top-5
STOP	100.0%	100.0%
SEQ	100.0%	100.0%

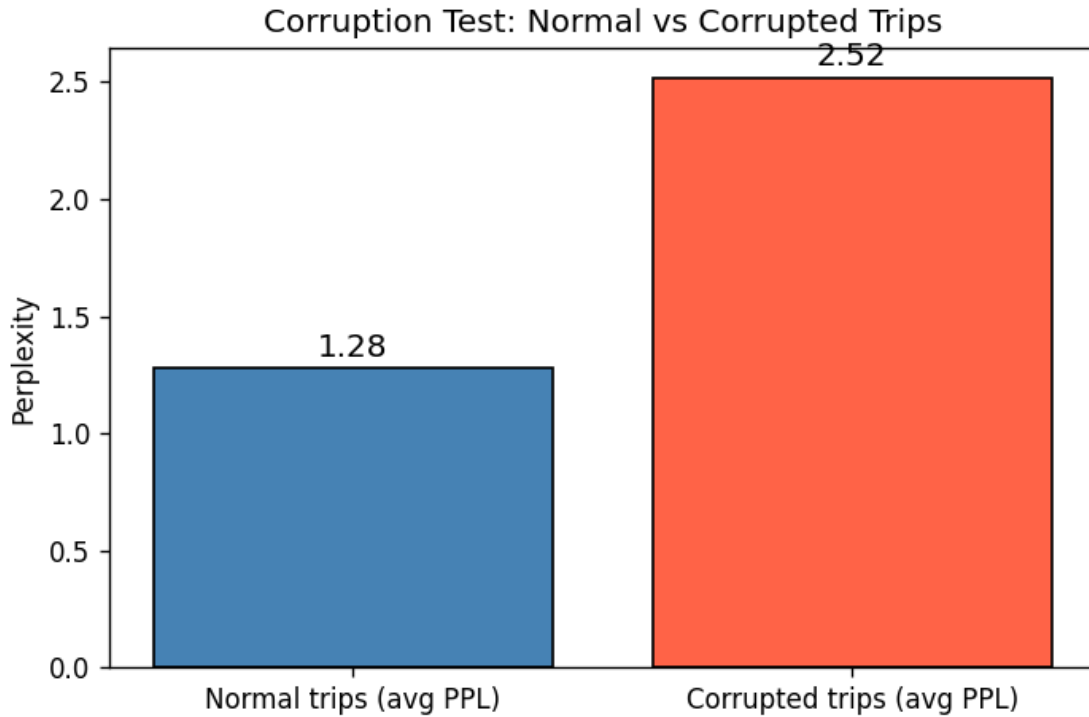
1.4.2 Corruption Test

To validate that the model learned network topology, real trips were corrupted by replacing middle stop tokens with random stops from other routes. If the model learned real structure, corrupted trips should have significantly higher perplexity.

```
[6]: corruption_data = {
      'Normal trips (avg PPL)': 1.28,
      'Corrupted trips (avg PPL)': 2.52,
    }

fig, ax = plt.subplots(figsize=(6, 4))
bars = ax.bar(corruption_data.keys(), corruption_data.values(),
              color=['steelblue', 'tomato'], edgecolor='black')
ax.set_ylabel('Perplexity')
ax.set_title('Corruption Test: Normal vs Corrupted Trips')
for bar, val in zip(bars, corruption_data.values()):
    ax.text(bar.get_x() + bar.get_width()/2, bar.get_height() + 0.05,
            f'{val:.2f}', ha='center', fontsize=12)
plt.tight_layout()
plt.show()

print("Corrupted trips show ~2x higher perplexity.")
print("The model detects out-of-place stops in the network.")
```



Corrupted trips show ~2x higher perplexity.
 The model detects out-of-place stops in the network.

1.5 4. Route Embeddings

The model's token embedding layer maps each route to a 384-dimensional vector. Routes that share similar stop patterns, timing, and network position should cluster together in embedding space.

```
[7]: import sys
      sys.path.insert(0, ROOT)
      import torch
      from model import GPTConfig, GPT
      from sklearn.decomposition import PCA

      ckpt = torch.load(os.path.join(ROOT, 'out-mtl-transit', 'ckpt.pt'),
        ↪map_location='cpu', weights_only=False)
      conf = GPTConfig(**ckpt['model_args'])
      model = GPT(conf)
      state_dict = ckpt['model']
      for k in list(state_dict.keys()):
        if k.startswith('_orig_mod.'):
          state_dict[k[len('_orig_mod.'):]] = state_dict.pop(k)
      model.load_state_dict(state_dict)
      model.eval()
```

```

wte = model.transformer.wte.weight.detach().numpy()
stoi = meta['stoi']

route_tokens = sorted([t for t in stoi if t.startswith('<ROUTE_')],
                       key=lambda t: int(t.split('_')[1].rstrip('>')))
route_ids = [stoi[t] for t in route_tokens]
route_vecs = wte[route_ids]
route_labels = [t.replace('<ROUTE_', '').rstrip('>') for t in route_tokens]

pca = PCA(n_components=2)
coords = pca.fit_transform(route_vecs)

fig, ax = plt.subplots(figsize=(14, 10))
ax.scatter(coords[:, 0], coords[:, 1], s=25, alpha=0.7, c='steelblue')

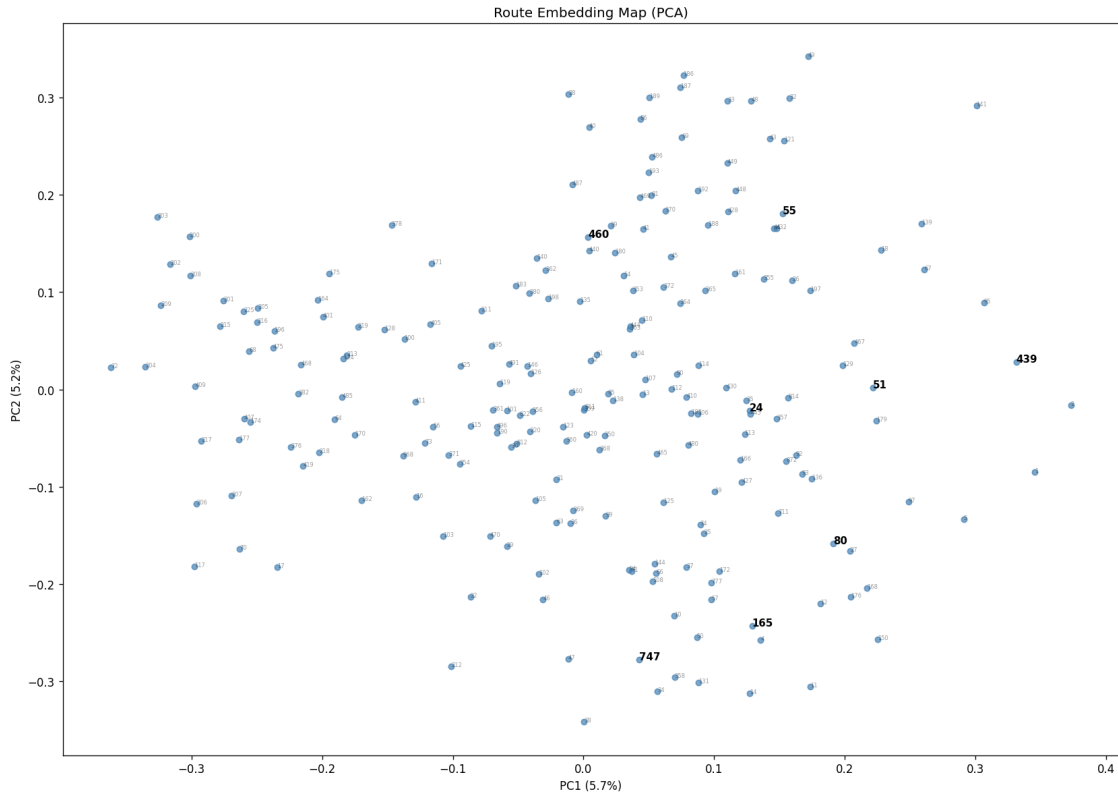
highlight = {'80', '24', '55', '747', '439', '165', '51', '460'}
for i, label in enumerate(route_labels):
    alpha = 1.0 if label in highlight else 0.4
    size = 9 if label in highlight else 5
    weight = 'bold' if label in highlight else 'normal'
    ax.annotate(label, (coords[i, 0], coords[i, 1]),
                fontsize=size, alpha=alpha, fontweight=weight)

ax.set_title('Route Embedding Map (PCA)')
ax.set_xlabel(f'PC1 ({pca.explained_variance_ratio_[0]:.1%})')
ax.set_ylabel(f'PC2 ({pca.explained_variance_ratio_[1]:.1%})')
plt.tight_layout()
plt.show()

print(f"{len(route_tokens)} routes projected. Highlighted: {'', '}'
      →join(sorted(highlight))}")

```

number of parameters: 14.71M



215 routes projected. Highlighted: 165, 24, 439, 460, 51, 55, 747, 80

1.6 5. Network Bottlenecks

Segments ranked by **coefficient of variation** (std / mean) of travel time across trips. High CV indicates inconsistent travel times — candidate operational bottlenecks.

```
[8]: bottlenecks = pd.read_csv(os.path.join(ANALYSIS_DIR, 'network_bottlenecks.csv'))
top_bn = bottlenecks.head(15).copy()
top_bn['label'] = ('R' + top_bn['route_id'].astype(str) + ': ' +
                  top_bn['stop_A_name'].str[:15] + ' -> ' +
                  top_bn['stop_B_name'].str[:15])

fig, axes = plt.subplots(1, 2, figsize=(16, 6))

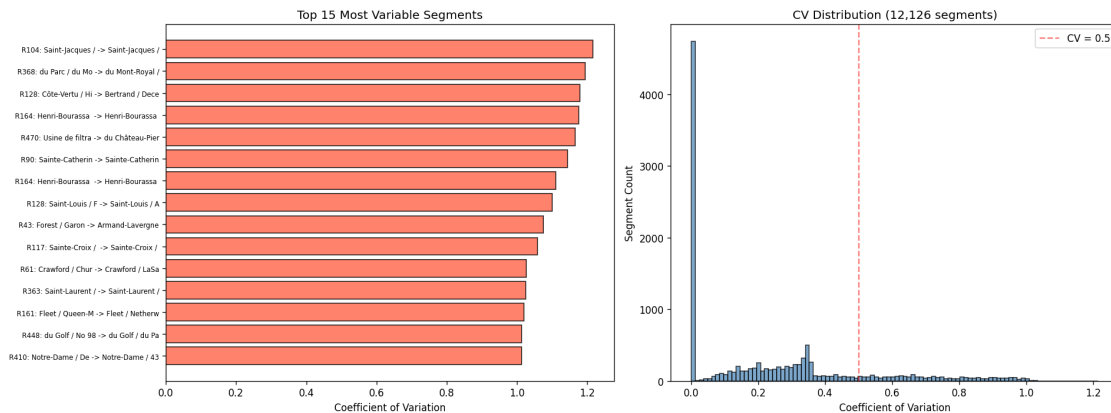
axes[0].barh(range(len(top_bn)), top_bn['cv'], color='tomato',
             edgecolor='black', alpha=0.8)
axes[0].set_yticks(range(len(top_bn)))
axes[0].set_yticklabels(top_bn['label'], fontsize=7)
axes[0].set_xlabel('Coefficient of Variation')
axes[0].set_title('Top 15 Most Variable Segments')
axes[0].invert_yaxis()
```

```

axes[1].hist(bottlenecks['cv'], bins=100, edgecolor='black', alpha=0.7,
↳color='steelblue')
axes[1].set_xlabel('Coefficient of Variation')
axes[1].set_ylabel('Segment Count')
axes[1].set_title(f'CV Distribution ({len(bottlenecks):,} segments)')
axes[1].axvline(x=0.5, color='red', linestyle='--', alpha=0.5, label='CV = 0.5')
axes[1].legend()

plt.tight_layout()
plt.show()

```



1.7 6. Recovery Zones (Schedule Padding)

Segments where scheduled travel time is significantly higher than neighboring segments. A high **padding ratio** (segment time / average of neighbors) indicates intentional schedule slack for delay recovery.

```

[9]: recovery = pd.read_csv(os.path.join(ANALYSIS_DIR, 'recovery_zones.csv'))
padded = recovery[recovery['padding_ratio'] >= 1.5].sort_values('padding_ratio',
↳ascending=False)

fig, axes = plt.subplots(1, 2, figsize=(14, 5))

axes[0].hist(recovery['padding_ratio'].clip(upper=10), bins=80,
edgecolor='black', alpha=0.7, color='steelblue')
axes[0].axvline(x=1.5, color='red', linestyle='--', label='Threshold (1.5)')
axes[0].set_xlabel('Padding Ratio')
axes[0].set_ylabel('Segment Count')
axes[0].set_title('Padding Ratio Distribution')
axes[0].legend()

```

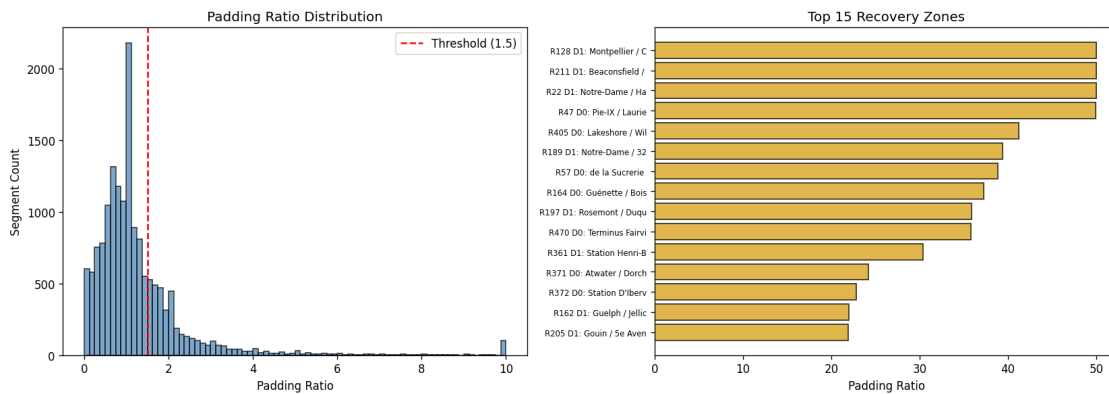
```

top_rz = padded.head(15).copy()
top_rz['label'] = ('R' + top_rz['route_id'].astype(str) + ' D' +
↳top_rz['direction_id'].astype(str) +
                ': ' + top_rz['stop_A_name'].str[:15])
top_rz = top_rz.iloc[::-1]
axes[1].barh(range(len(top_rz)), top_rz['padding_ratio'].clip(upper=50),
             color='goldenrod', edgecolor='black', alpha=0.8)
axes[1].set_yticks(range(len(top_rz)))
axes[1].set_yticklabels(top_rz['label'], fontsize=7)
axes[1].set_xlabel('Padding Ratio')
axes[1].set_title('Top 15 Recovery Zones')

plt.tight_layout()
plt.show()

print(f"Total segments: {len(recovery):,}")
print(f"Segments with ratio >= 1.5: {len(padded):,}")

```



Total segments: 15,745
Segments with ratio >= 1.5: 3,986

1.8 7. Route Variants

Many STM routes have multiple stop-sequence variants: short turns, branch termini, and alternate routing. GTFS encodes these as different trip patterns under the same route ID.

```

[12]: variants = pd.read_csv(os.path.join(ANALYSIS_DIR, 'route_variants.csv'))
variants_per_rd = variants.groupby(['route_id', 'direction_id']).size().
↳reset_index(name='n_variants')
multi = variants_per_rd[variants_per_rd['n_variants'] > 1]

fig, axes = plt.subplots(1, 2, figsize=(14, 5))

```

```

axes[0].hist(variants_per_rd['n_variants'], bins=range(1,
↳variants_per_rd['n_variants'].max() + 2),
            edgecolor='black', alpha=0.7, color='steelblue', align='left')
axes[0].set_xlabel('Number of Variants')
axes[0].set_ylabel('Route-Direction Count')
axes[0].set_title('Variants per Route-Direction')

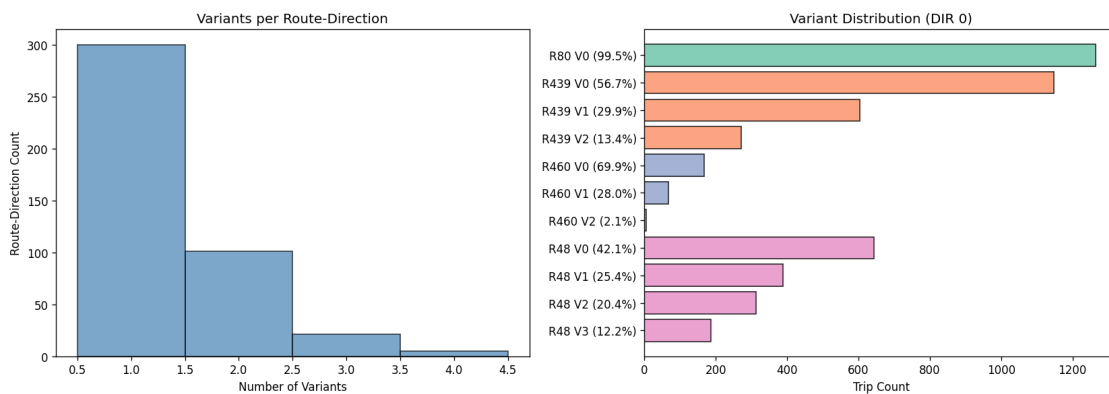
example_routes = [80, 439, 460, 48]
for ax_idx, rid in enumerate(example_routes[:4]):
    rv = variants[variants['route_id'] == rid].sort_values(['direction_id',
↳'trip_count'], ascending=[True, False])
    if not rv.empty:
        d0 = rv[rv['direction_id'] == 0]
        if not d0.empty:
            labels = [f"V{r['variant_id']} ({r['pct']}%)" for _, r in d0.
↳iterrows()]
            sizes = d0['trip_count'].values
            axes[1].barh([f'R{rid} ' + l for l in labels], sizes,
                        color=plt.cm.Set2(ax_idx), edgecolor='black', alpha=0.8)

axes[1].set_xlabel('Trip Count')
axes[1].set_title('Variant Distribution (DIR 0)')
axes[1].invert_yaxis()

plt.tight_layout()
plt.show()

print(f"Total route-directions: {len(variants_per_rd)}")
print(f"With multiple variants: {len(multi)}")

```



Total route-directions: 427
With multiple variants: 127

1.9 8. Route Timing Signatures

Travel time profile along Route 80, showing mean inter-stop time and variance at each segment. Spikes indicate segments where buses take longer or timing is less predictable.

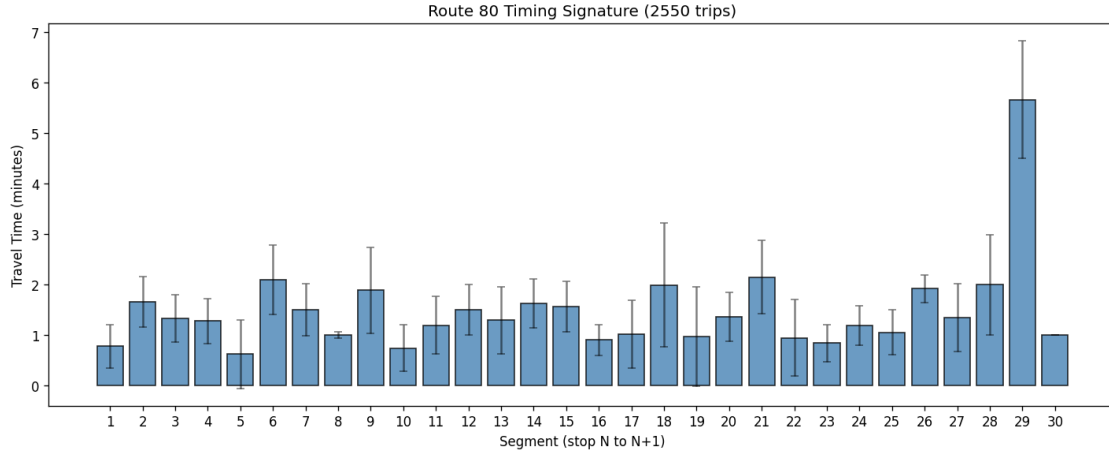
```
[13]: from collections import defaultdict

route_deltas = defaultdict(list)
for line in lines:
    tokens = line.strip().split()
    route = None
    times = []
    for tok in tokens:
        if tok.startswith('<ROUTE_'):
            route = tok.split('_')[1].rstrip('>')
        elif tok.startswith('<T_'):
            times.append(int(tok.split('_')[1].rstrip('>')))
    if route == '80' and len(times) >= 2:
        route_deltas['80'].append([times[i+1] - times[i] for i in
        ↪range(len(times)-1)])

all_deltas = route_deltas['80']
max_stops = max(len(d) for d in all_deltas)

means, stds = [], []
for pos in range(min(max_stops, 30)):
    vals = [d[pos] for d in all_deltas if pos < len(d)]
    means.append(np.mean(vals))
    stds.append(np.std(vals))

positions = range(1, len(means) + 1)
fig, ax = plt.subplots(figsize=(12, 5))
ax.bar(positions, means, yerr=stds, capsize=3, color='steelblue',
        edgcolor='black', alpha=0.8, error_kw={'alpha': 0.5})
ax.set_xlabel('Segment (stop N to N+1)')
ax.set_ylabel('Travel Time (minutes)')
ax.set_title(f'Route 80 Timing Signature ({len(all_deltas)} trips)')
ax.set_xticks(positions)
plt.tight_layout()
plt.show()
```



1.10 9. Conclusions

Trip grammar exists in transit networks. STM trips follow predictable patterns of route, direction, stop ordering, and timing that a transformer can learn with high accuracy.

The model learned route topology. 98.8% next-stop prediction accuracy and corruption detection demonstrate genuine understanding of which stops follow which on each route.

Routes cluster into behavioral families. PCA of route embeddings reveals that the model groups routes by shared structural properties — not just geographic proximity.

Bottlenecks and recovery zones are detectable from schedule data alone. Coefficient of variation analysis identifies the most operationally unstable segments. Padding ratio analysis reveals where STM intentionally builds schedule slack.

Route variants encode hidden operational patterns. Many routes have short-turn services, branch termini, and alternate paths that only emerge from systematic trip-sequence analysis.

GTFS static data alone reveals structural patterns that would typically require real-time operational data or manual inspection to discover. Phase 2 will integrate GTFS-RT data to compare scheduled vs actual performance.