# On Simulation and Optimization of Freeway Network Operations

**B. Wiwatanapataphee, Yong Hong Wu, C Gu**
**Curtin University of Technology**

# Progress from last PSG meeting

1  Freeway traffic control via Ramp Metering and Variable Speed Limit using a mesoscopic model

2.  Control of ramp metering based on reinforcement learning

# Freeway Traffic Control via Ramp Metering and Variable Speed Limit using a mesoscopic model

**B. Wiwatanapataphee & YH Wu**
**Curtin University of Technology**

# Background

**PERTH**

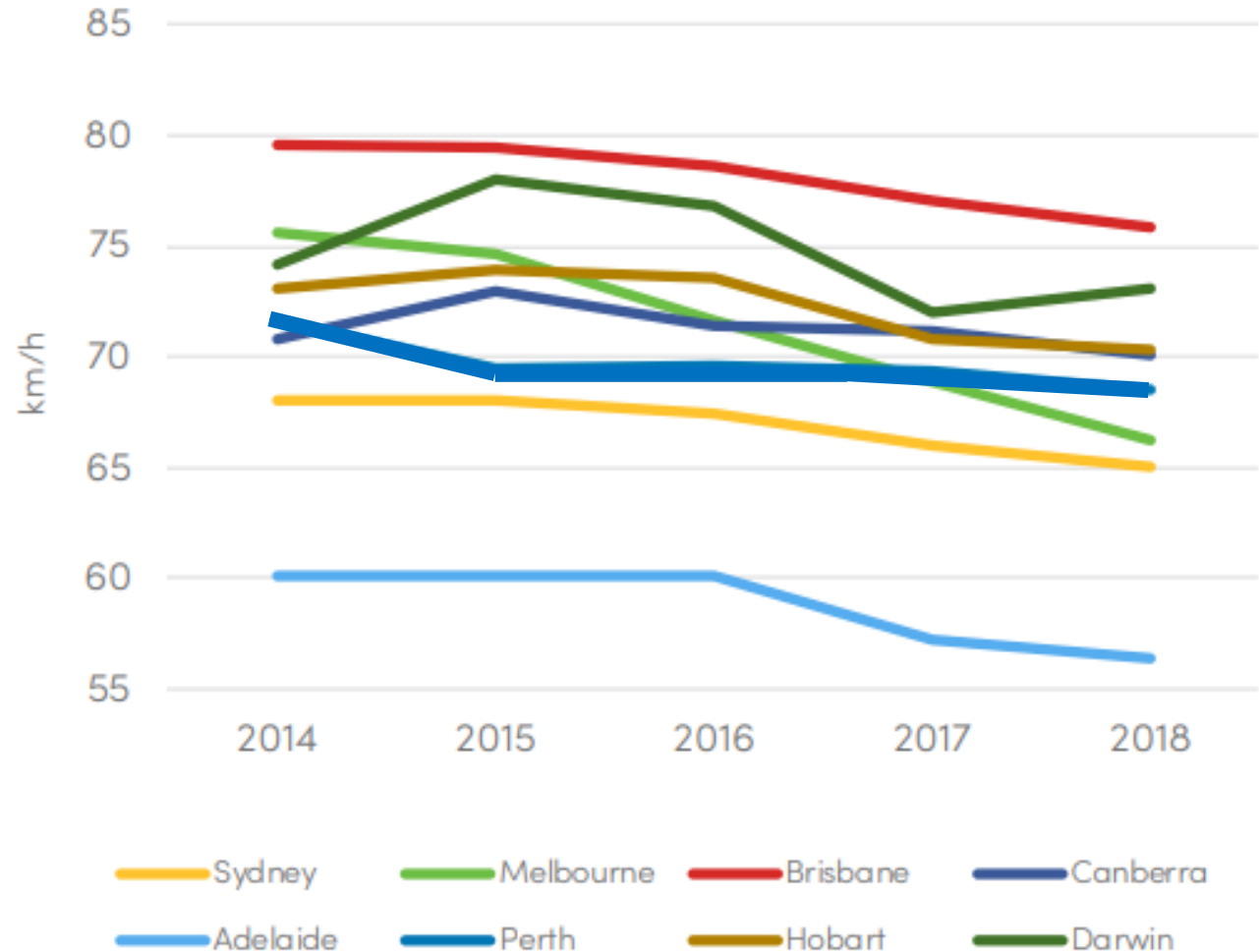| | | |
|---|---|---|
| 🎚️ | AVG. SPEED | **61.6** KM/H |
| 🚗 | CONGESTED SPEED (% OF FREE FLOW) | **94.8%** |
| 🕐 | VARIABILITY | **24.3%** |



Average free flow speeds

https://www.mainroads.wa.gov.au/about-main-roads/news-media/smart-freeway-technology-upgrades/

# Traffic monitoring and Control systems

❖ Traffic congestion has a significant impact on economic activity throughout much of the world.

❖ An essential step towards active congestion control is the creation of accurate, reliable traffic monitoring and control systems.

❖ These systems usually run algorithms which rely on mathematical models of traffic used to power estimation and control schemes.
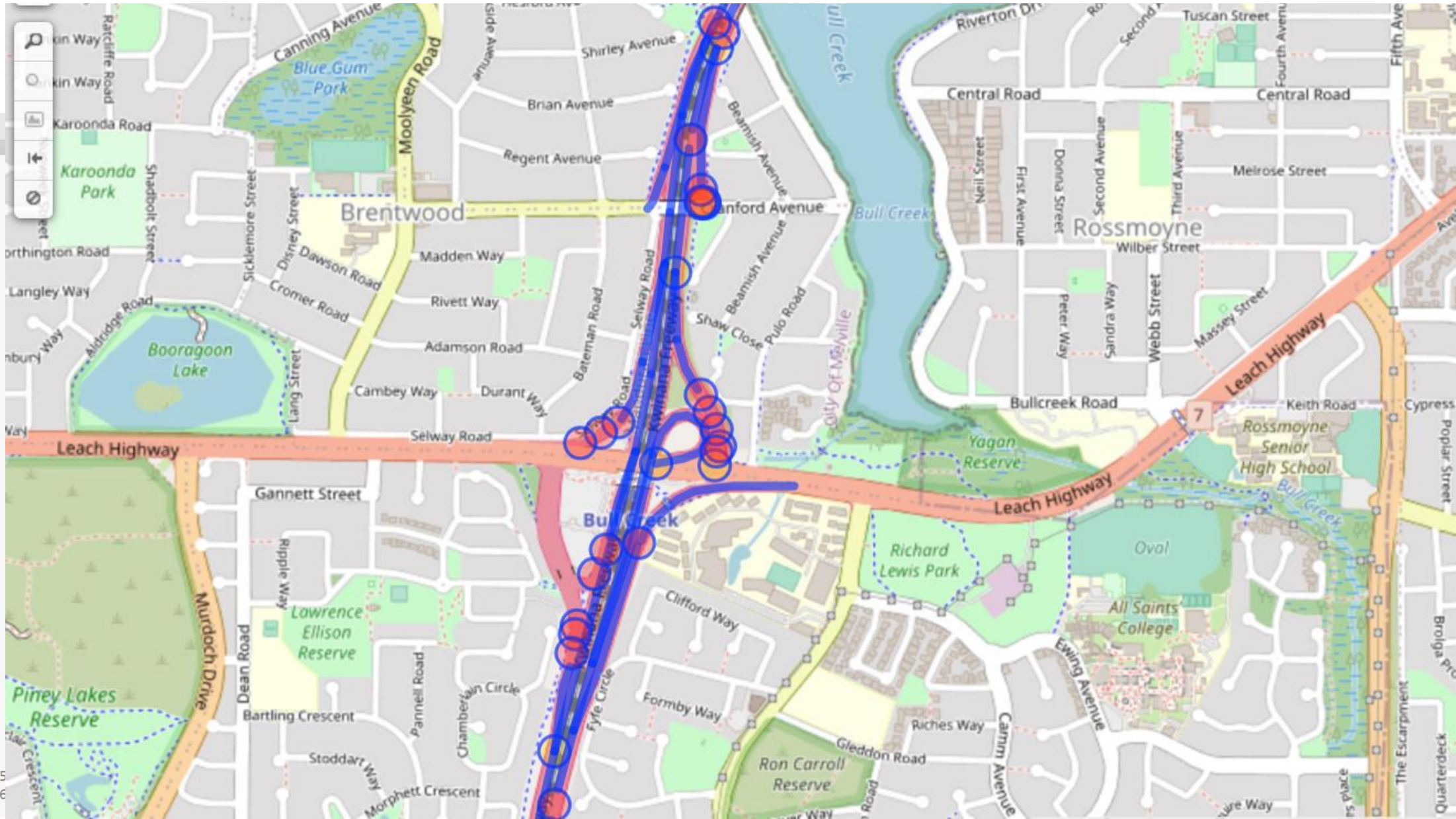
# Road Network

10

153 cells

Incoming flow
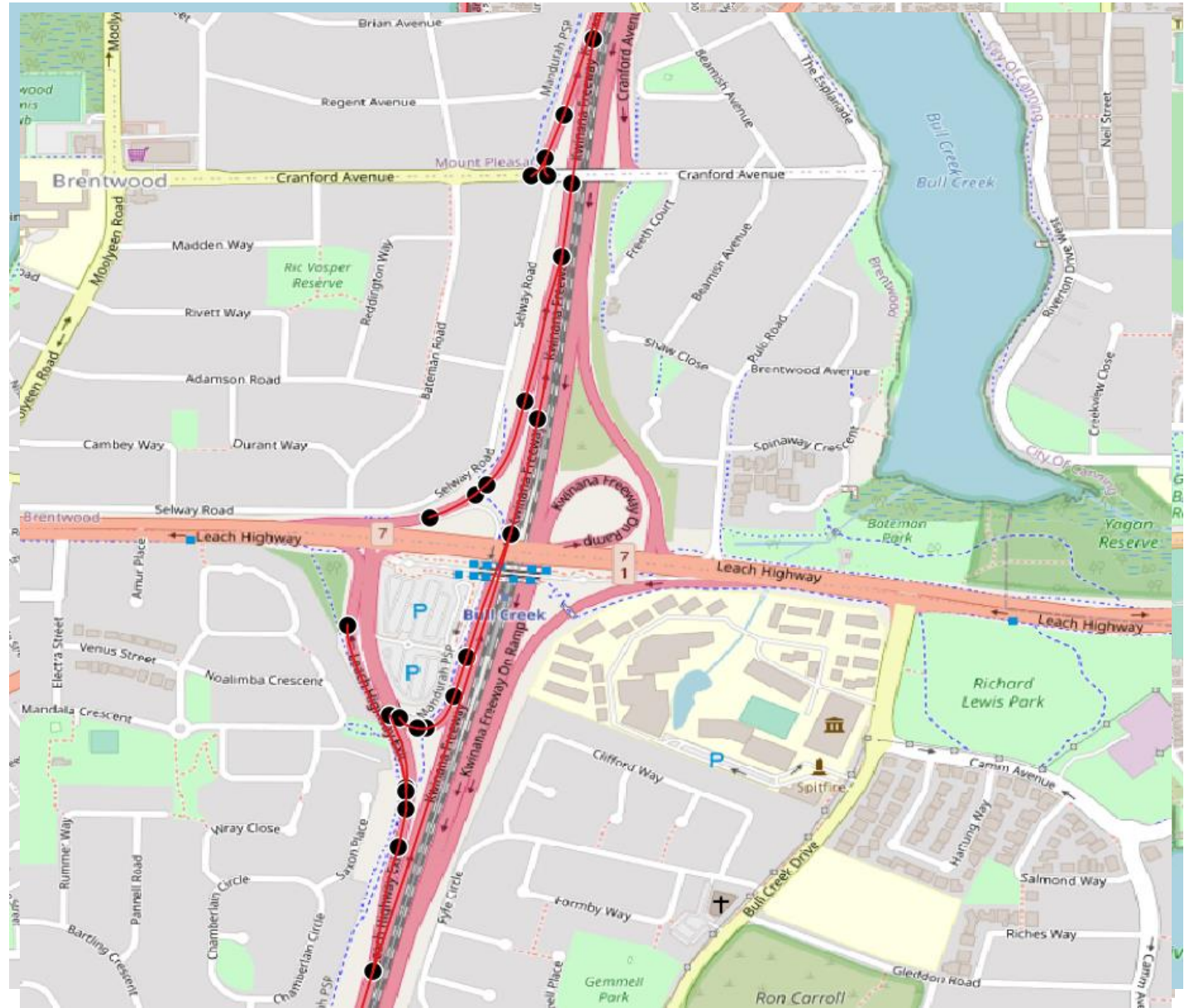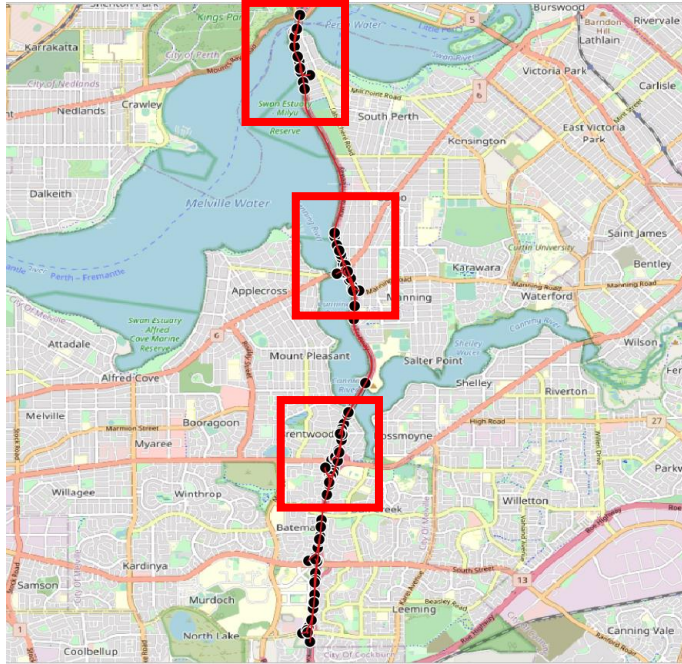
Outgoing flow

① ②

**Legend**

| Source | Ordinary | Diverging | Merging | Sink |

**cells**

| Ordinary | Merging | Diverging |

**Links**
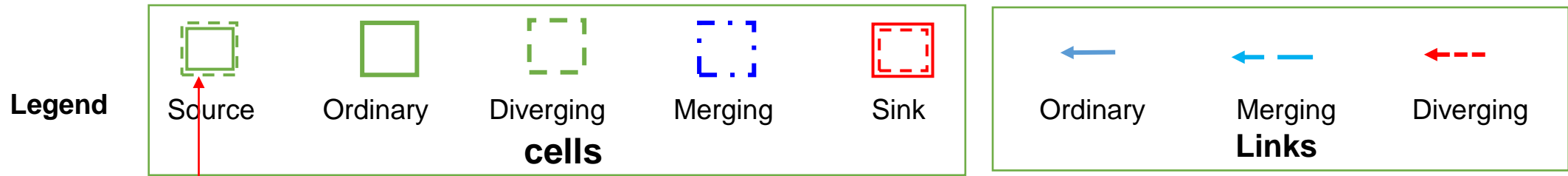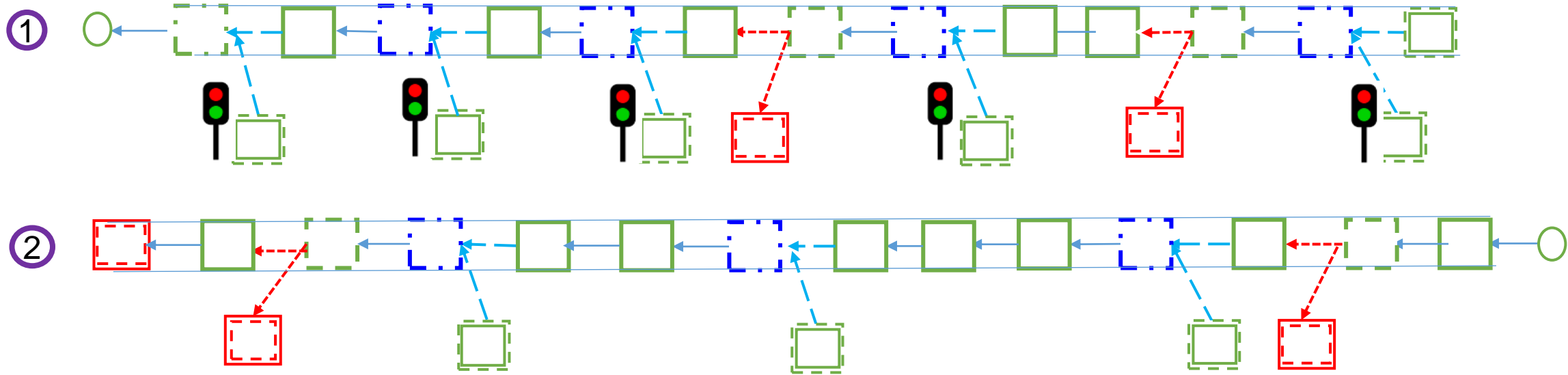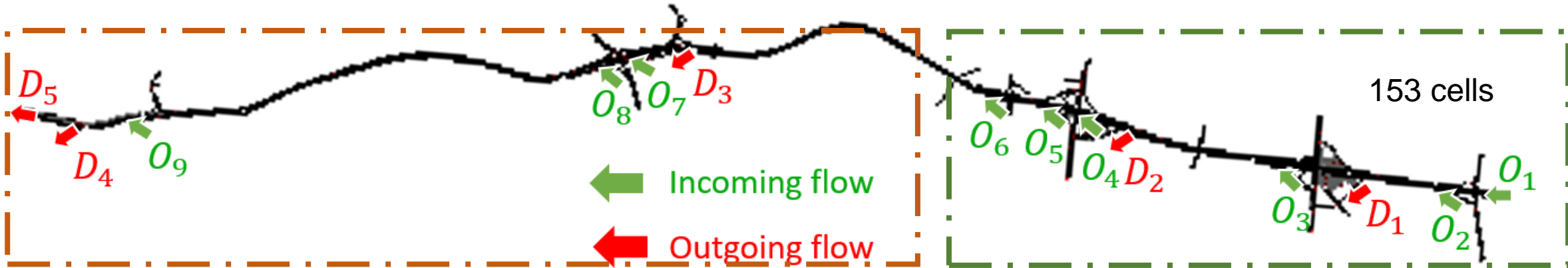
Estimated demand using CNN deep learning every 15 minutes

11

# The DE-CTM optimisation model:

$$TTS = \min \left[ \sum_{t=\tau}^{T} \sum_{i \in C} TL_i k_i(\alpha, \beta, t) \right]$$

$\boldsymbol{\alpha} \in [\boldsymbol{\alpha_{min}}, \boldsymbol{\alpha_{max}}]$ : probability from upstream normal cell;

$\boldsymbol{1 - \alpha}$: probability from upstream merge cell

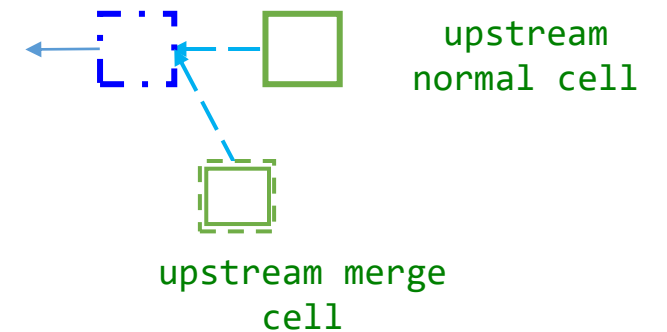$\boldsymbol{\beta} \in [\boldsymbol{\beta_{min}}, \boldsymbol{\beta_{max}}]$ : VSL rate

PDE:
$$\frac{n_i \partial k_i}{\partial t} + \frac{\partial(n_i q_i)}{\partial x} = g(\alpha, \beta, t)$$

$$\boldsymbol{r_{i,j}(t)} = \mathbf{min} \left[ (1 - \alpha)\boldsymbol{q_{i-1}(t)}; \quad \boldsymbol{Q_r}; \quad \boldsymbol{m_i(t-1)} + \frac{T}{L_j}\left(\boldsymbol{d_{i,j}} + \boldsymbol{r_{i,j}(t-1)}\right) \right], \boldsymbol{j} = 1, \dots, 8$$

Constraints:
$$\sum_{t=\tau}^{\tau+H_c-1} T \left\{ \sum_{i \in C_M} r_{i,j}(t) + n_i(q_i(t) - Q_i) \right\} \le 0, \qquad j = 1, \dots, 8$$

$$\sum_{t=\tau}^{\tau+H_c-1} T \left\{ \sum_{i \in C_{VSL}} n_i k_i(t) \left(\boldsymbol{\beta(i, t)} v_f\right) - n_i Q_i \right\} \le 0$$
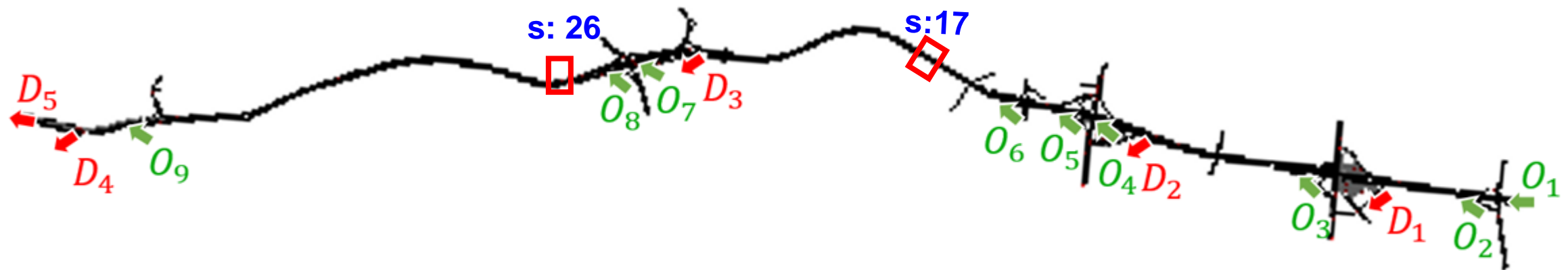
$$0 \le k_i(t) \le k_{jam}$$



upstream normal cell

upstream merge cell

# Some results $\alpha_m \in [-0.25, 0.25], \beta_{VSL} \in [0.25, 1.0]$

Let $q_{max}$ = 1800 vph on freeway    $q_{max}$= 1300 vph on on/off-ramp

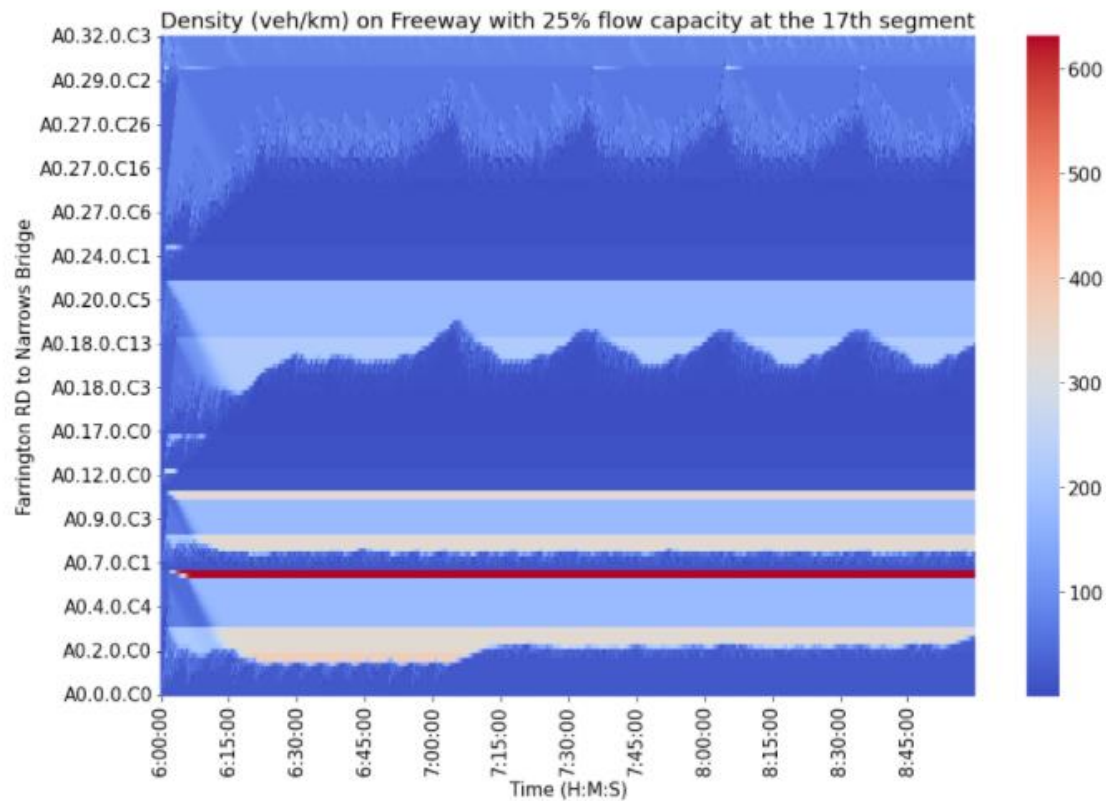| Case study | | Segment | Free speed (km/hr) | Capacity(veh/hr) |
|---|---|---|---|---|
| I | a | All segments | 100 | 1800 |
| | b | All segments | 70 | 1800 |
| II | a | 17 | 100 | 450 |
| | b | 26 | 100 | 450 |
| III | a | 17 | 70 | 450 |
| | b | 26 | 70 | 450 |

# Case I(b): Heatmap plot of density and flow rate

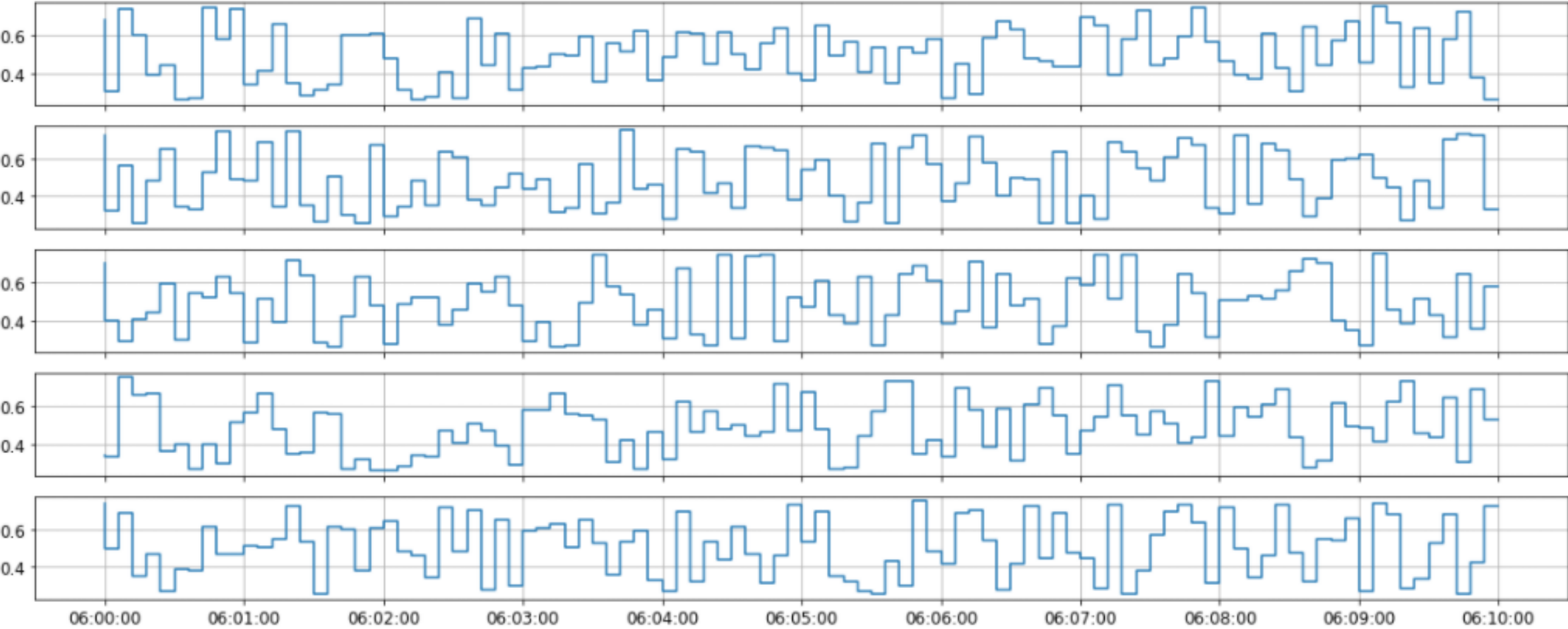**Free speed 100**

Traffic density (veh/km)                    Traffic flow rate (veh/hr)
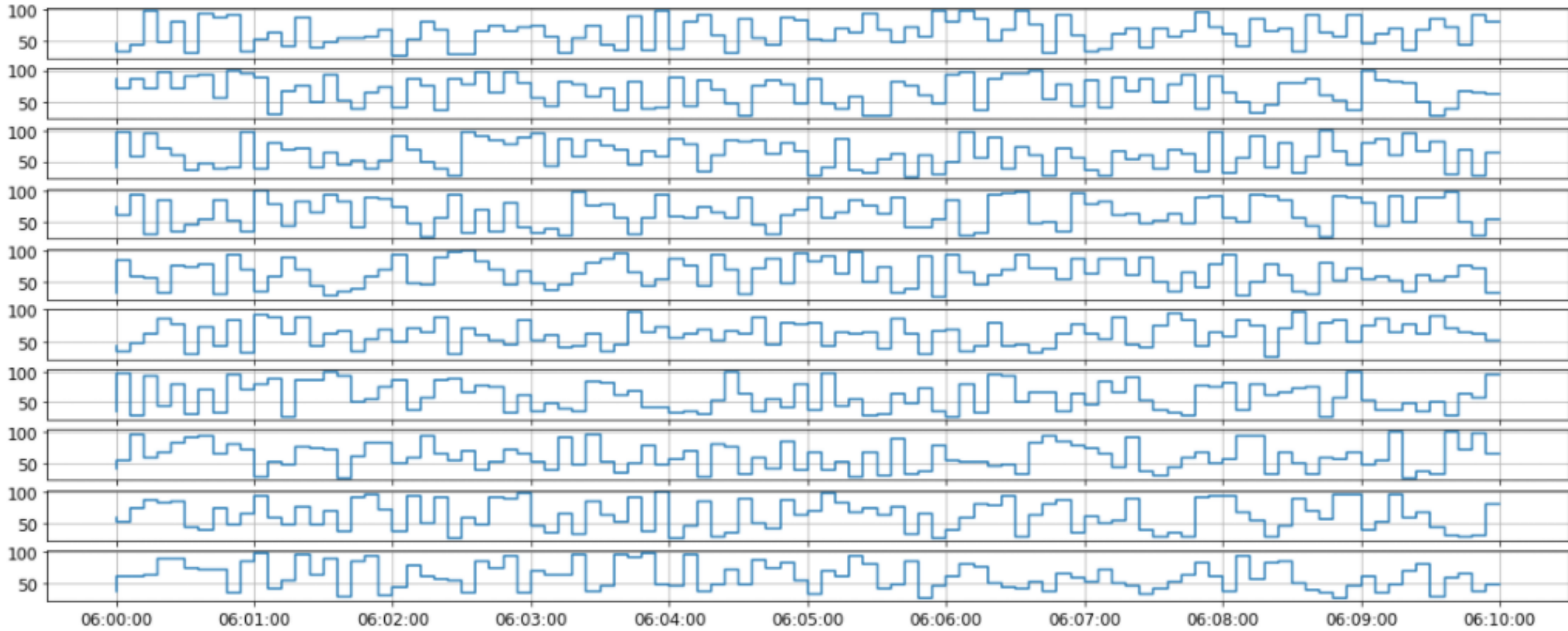
# Case I(a): RM control  and Variable speed limit

Ramp Metering at five on-ramps from Farrington to Cranford

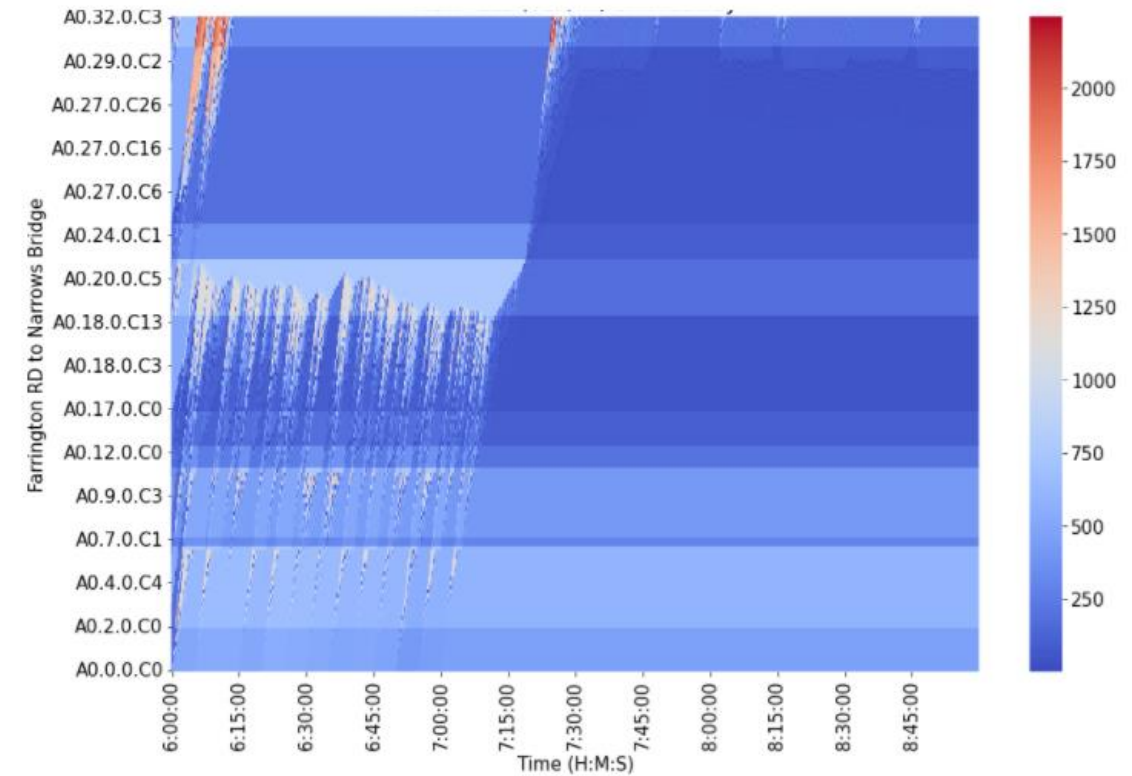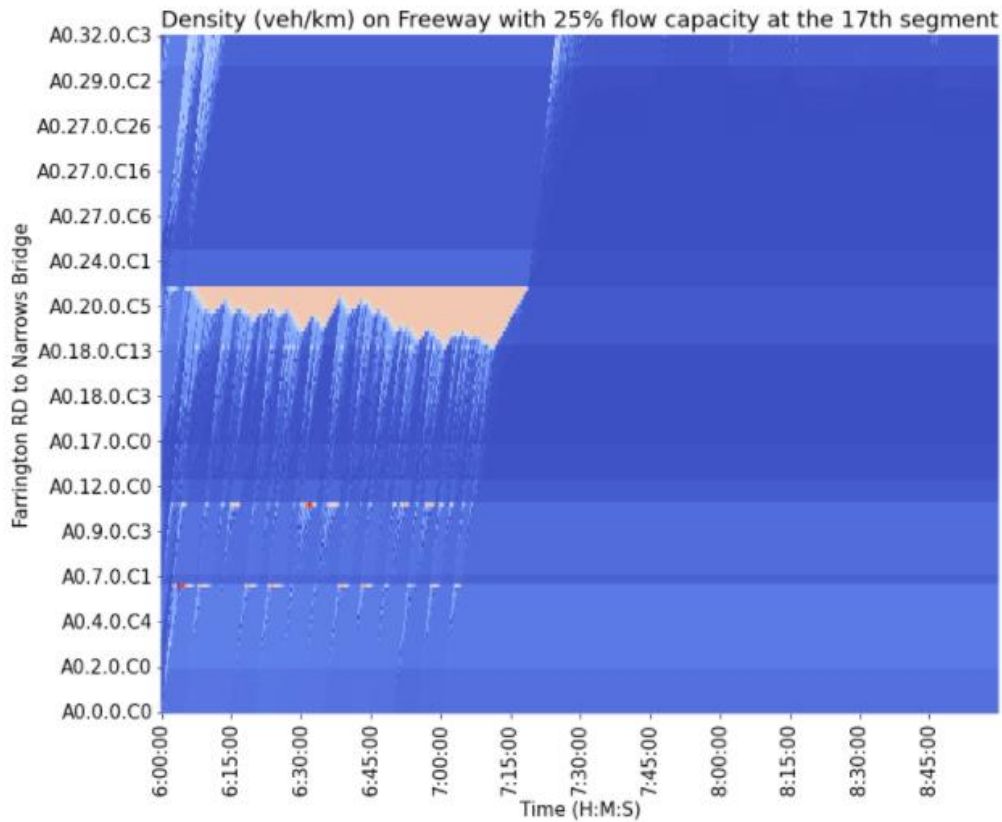Variable Speed Limits from Cranford to Mill Points

# Case I (a): Heatmap plot of density and flow rate
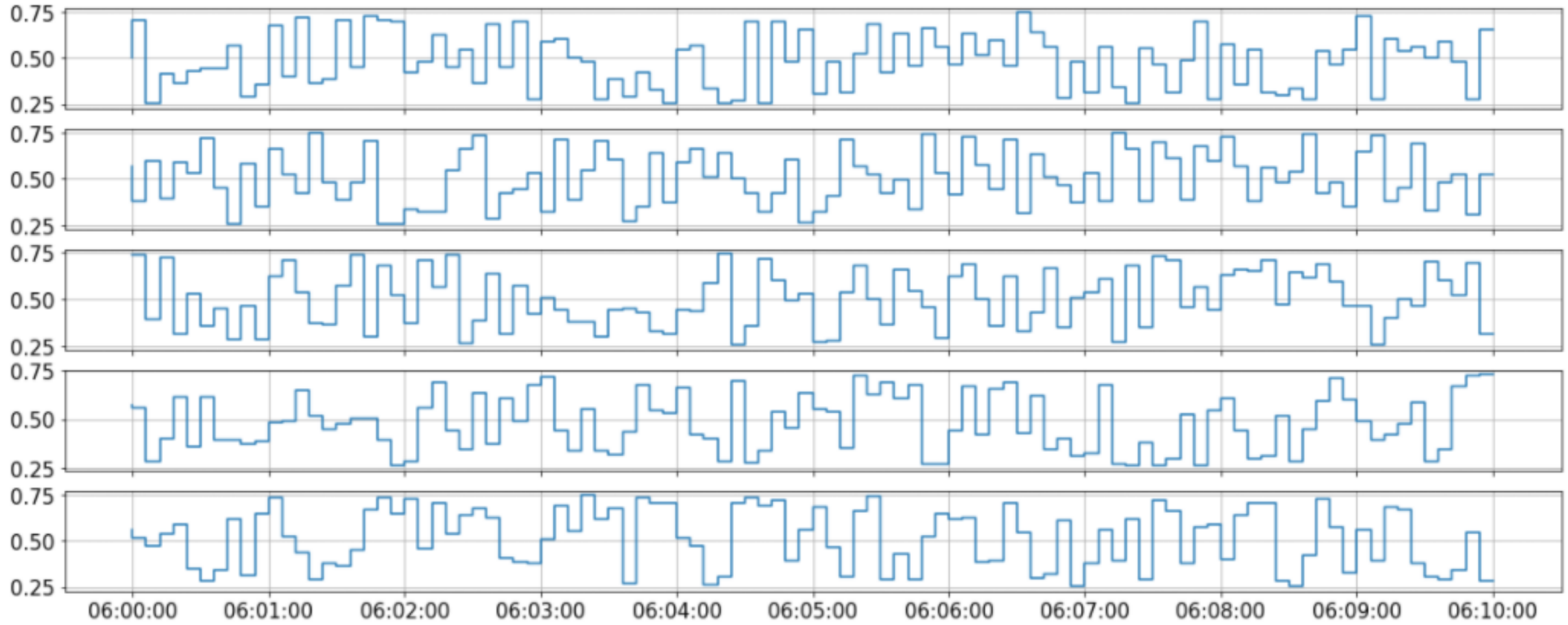
**Free speed 70**

Traffic density (veh/km)                    Traffic flow rate (veh/hr)
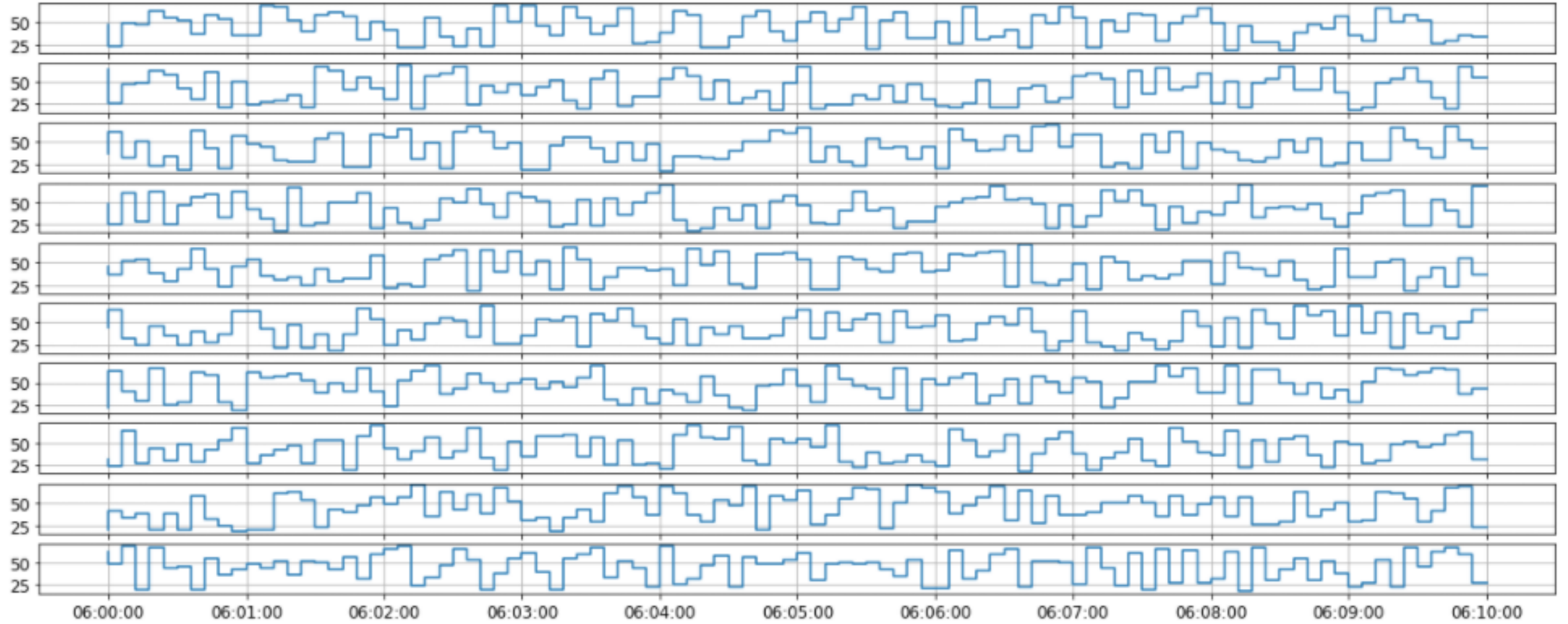
# Case I (b): RM control  and Variable speed limit

Ramp Metering at five on-ramps from Farrington to Cranford

Variable Speed Limits from Cranford to Mill Points

# Case II(a): Heatmap plot of density and flow rate

## Segment 17

Traffic density (veh/km)

Traffic flow rate (veh/hr)

# Case II(a): RM control and Variable speed limit

Ramp Metering at five on-ramps from Farrington to Cranford
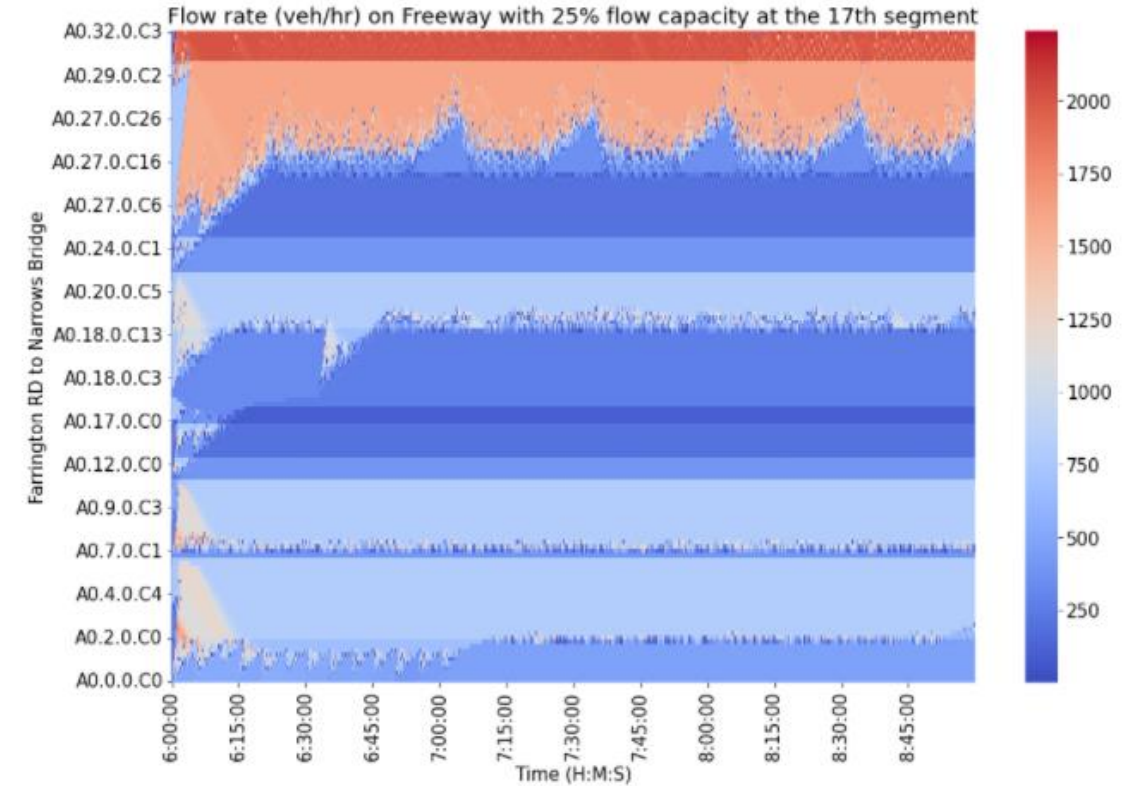
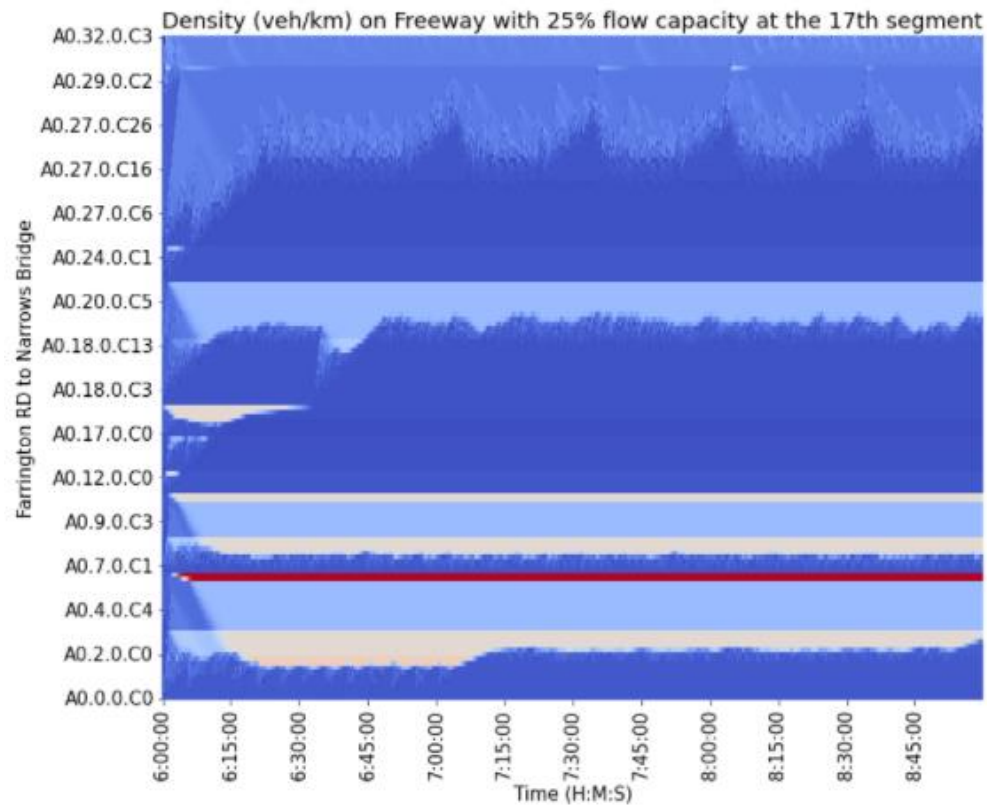Variable Speed Limits from Cranford to Mill Points

# Case II(b): Heatmap plot of density and flow rate
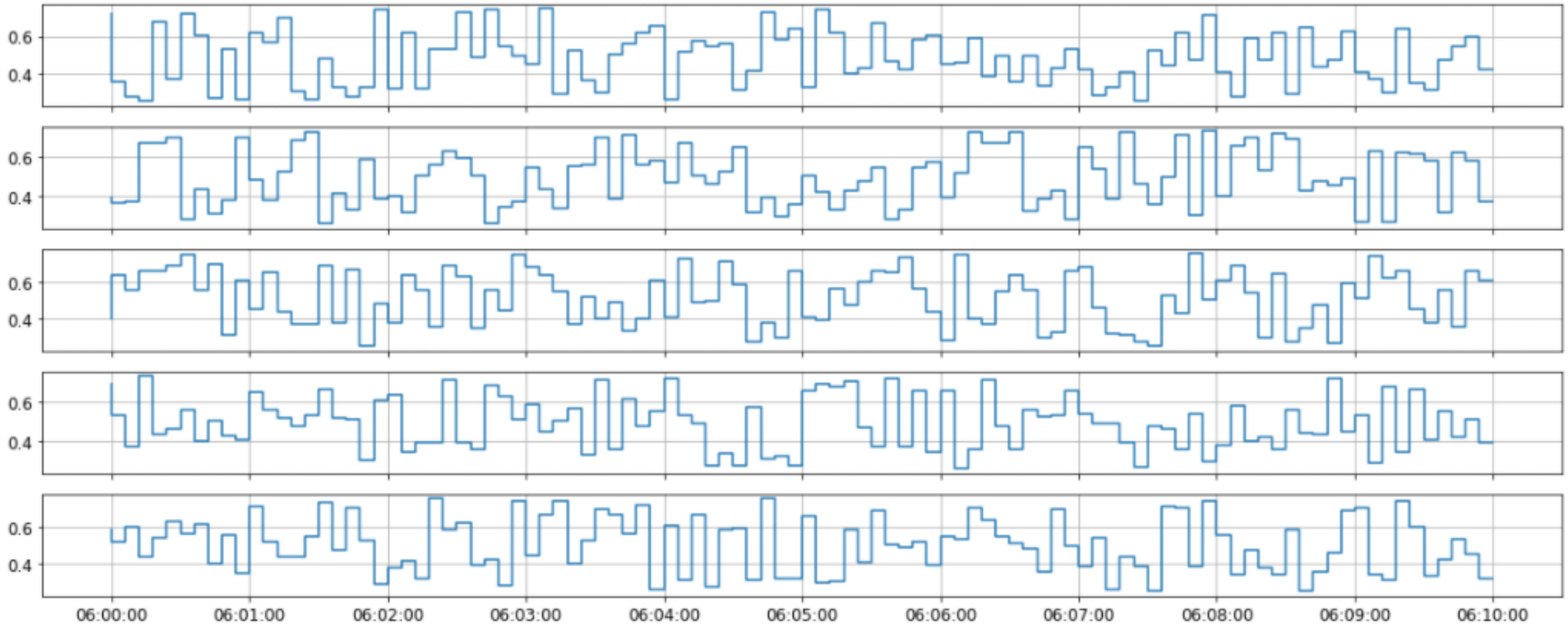
**Segment 26**

Traffic density (veh/km)
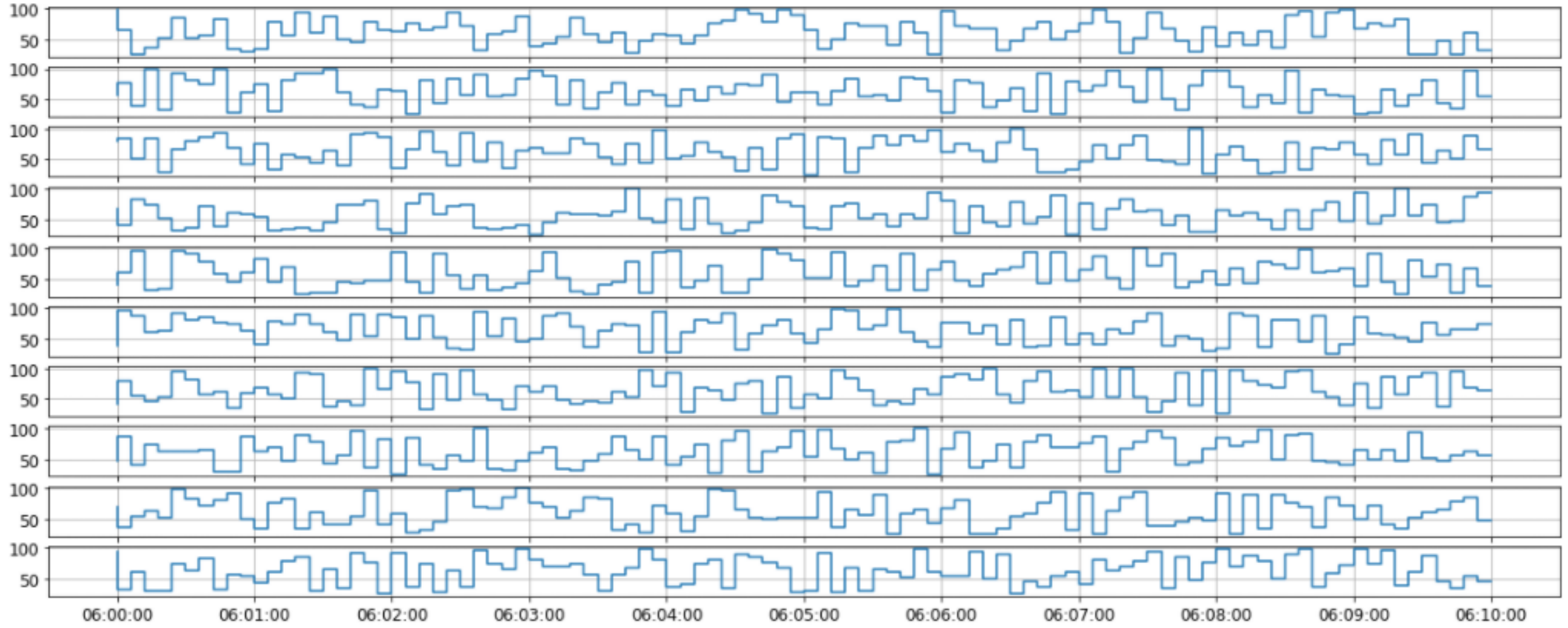
Traffic flow rate (veh/hr)

# Case II(b): RM control  and Variable speed limit

Ramp Metering at five on-ramps from Farrington to Cranford

Variable Speed Limits from Cranford to Mill Points

# Case III(a): Heatmap plot of density and flow rate
## Segment 17

Traffic density (veh/km)                              Traffic flow rate (veh/hr)

# Case III(a): RM control  and Variable speed limit
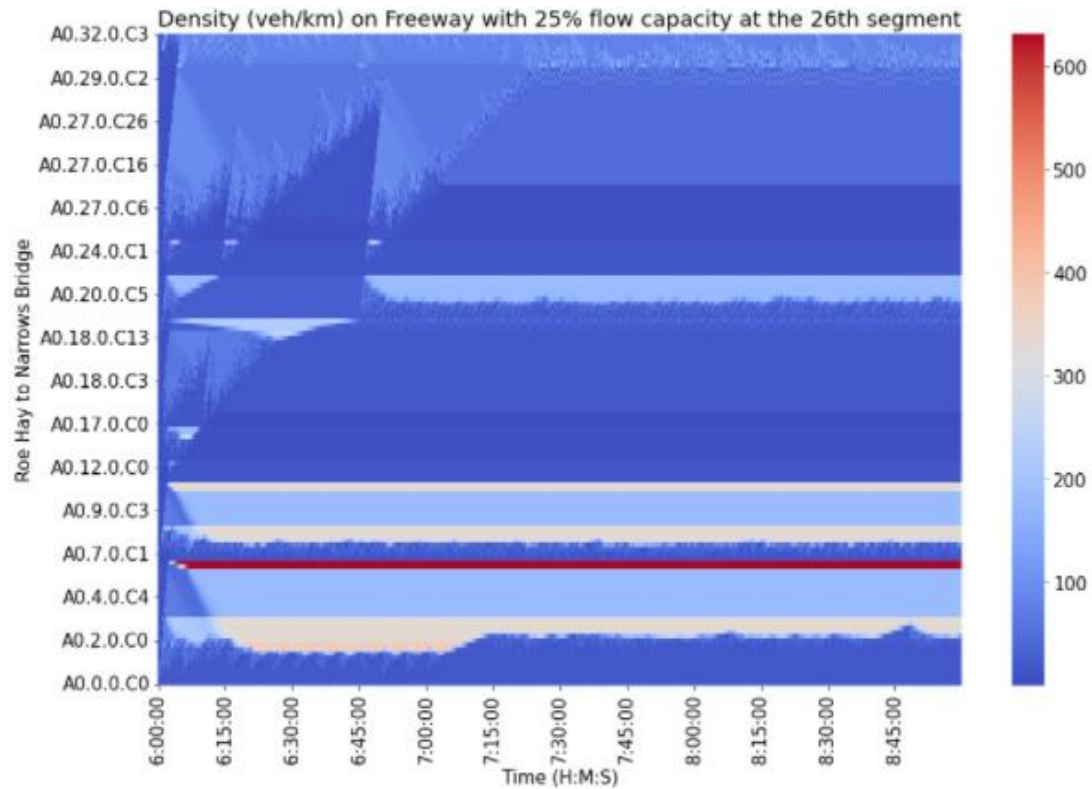
Ramp Metering at five on-ramps from Farrington to Cranford
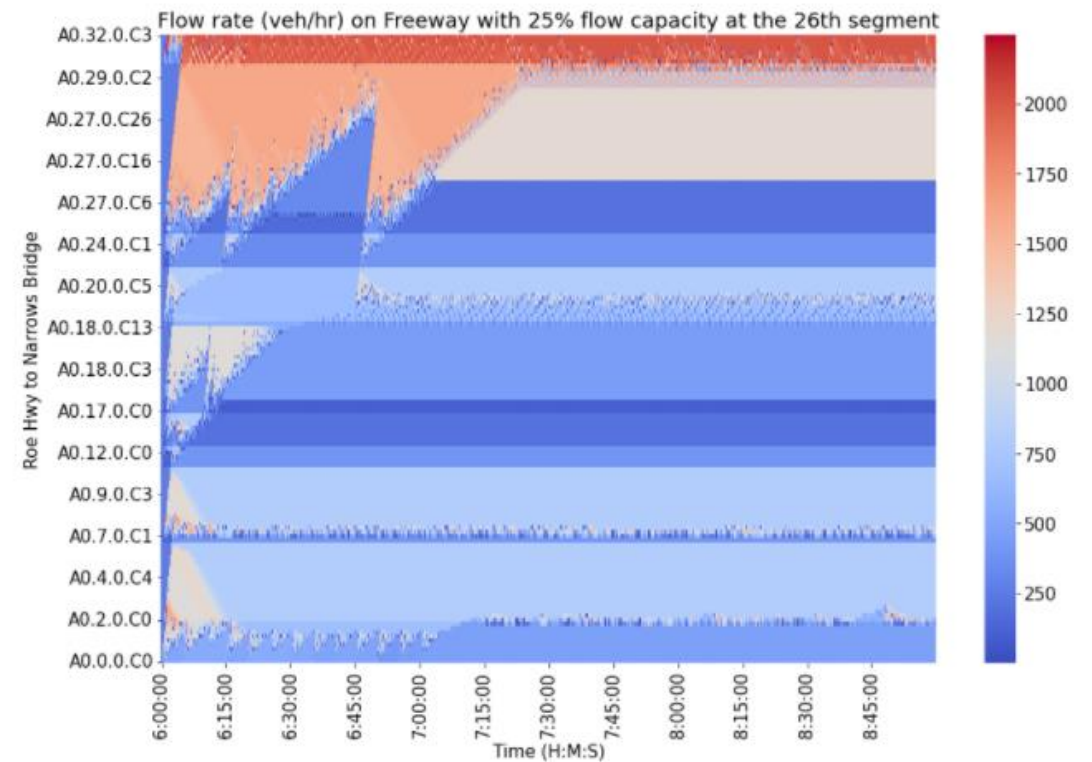
Variable Speed Limits from Cranford to Mill Points

# Case III(b): Heatmap plot of density and flow rate

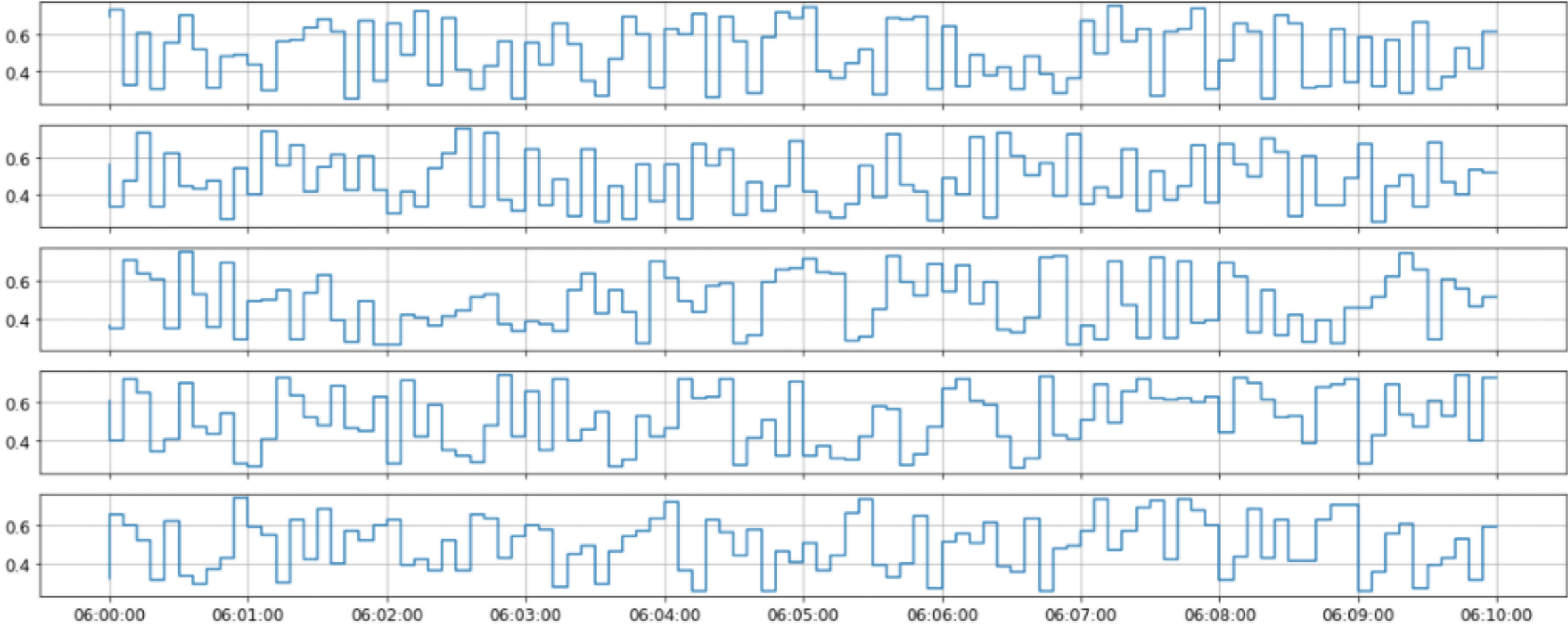**Segment 26**

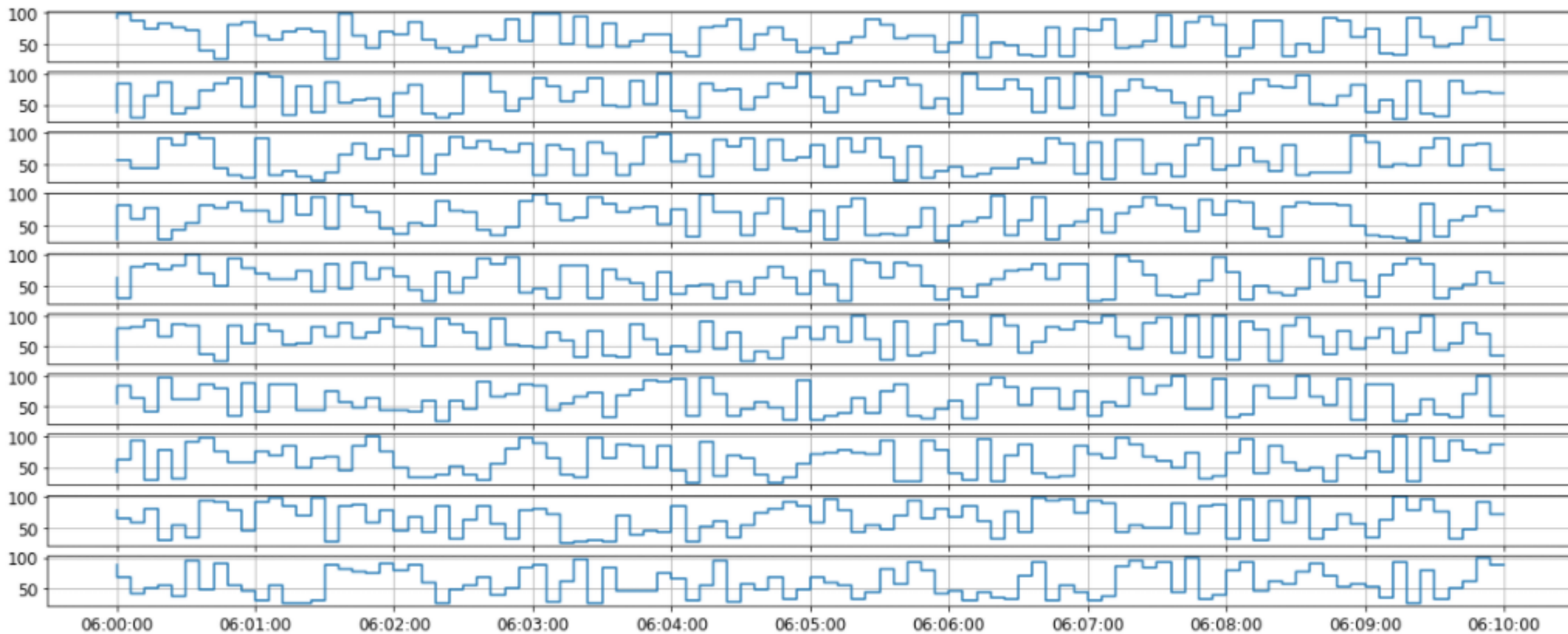Traffic density (veh/km)

Traffic flow rate (veh/hr)

# Case III(b): RM control  and Variable speed limit

# Ramp Metering at five on-ramps from Farrington to Cranford

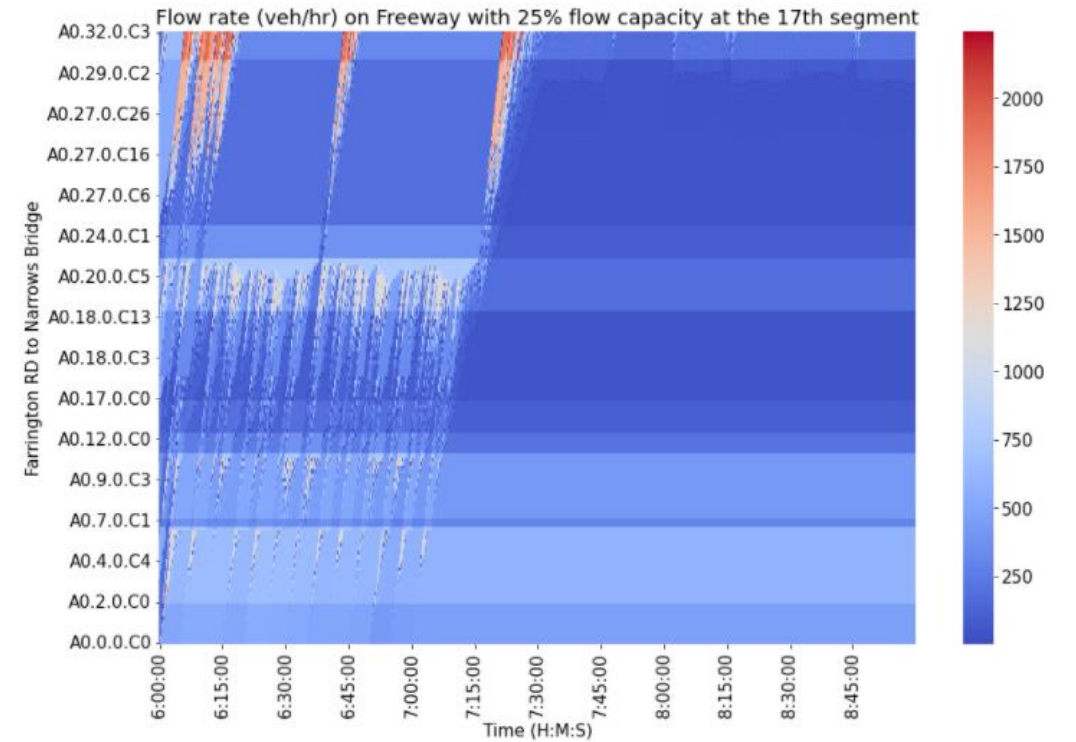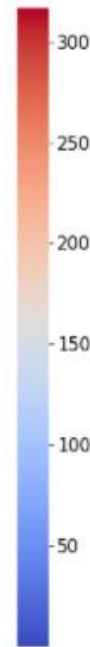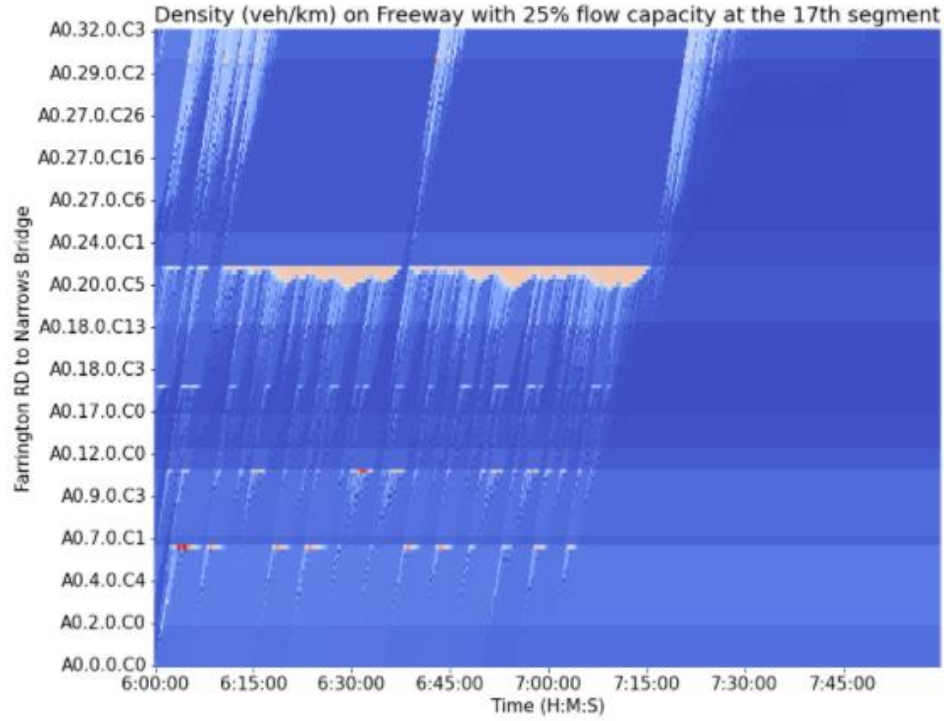Variable Speed Limits from Cranford to Mill Points

# 2. Control of ramp metering based on reinforcement learning

PhD Student: C. Gu;  Supervisors: YH Wu  & B Wiwatanapataphee

# Control of ramp metering based on reinforcement learning

state： queue length, mean
waiting time, mean
speed....

action： switching signal
phase, $a \in \{0,1\}$

policy: $\pi(a \mid s)$

$\pi(1 \mid s) = 0.8$

$\pi(0 \mid s) = 0.2$

reward: mean waiting time, mean time, total time spent....

state transition: old state $\longrightarrow$ new state $\quad S' \sim p(\cdot \,|\, s, a)$

- (state, action, reward) trajectory:

$$s_1, a_1, r_1, \quad s_2, a_2, r_2, \quad \cdots \quad, \quad s_n, a_n, r_n.$$

- One episode is from the the beginning to the end

$$s_1 \longrightarrow a_1 \longrightarrow s_2 \longrightarrow a_2 \longrightarrow s_3 \longrightarrow a_3 \longrightarrow s_4 \longrightarrow \cdots$$
$$a_1 \searrow r_1 \qquad a_2 \searrow r_2 \qquad a_3 \searrow r_3$$

# Randomness in Returns

**Definition:** Discounted return (at time $t$).

- $U_t = R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \cdots + \gamma^{n-t} R_n$.

At time $t$, the rewards, $R_t, \cdots, R_n$, are random, so the return $U_t$ is random.

- Reward $R_i$ depends on $S_i$ and $A_i$.

- States can be random: $\quad S_i \sim p(\,\cdot \mid s_{i-1},\ a_{i-1}\,)$.

- Actions can be random: $\quad A_i \sim \pi(\,\cdot \mid s_i\,)$.

- If either $S_i$ or $A_i$ is random, then $R_i$ is random.

# Action-Value Function $Q_\pi(s, a)$

**Definition:** Discounted return.

- $U_t = R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \cdots + \gamma^{n-t} R_n$.

**Definition:** Action-value function.

- $Q_\pi(s_t, a_t) = \mathbb{E}\left[\, U_t \mid S_t = s_t, A_t = a_t \,\right]$.

- $Q_\pi(s_t, a_t)$ depends on $s_t$, $a_t$, $\pi$, and $p$.
- $Q_\pi(s_t, a_t)$ is dependent of $S_{t+1}, \cdots, S_n$ and $A_{t+1}, \cdots, A_n$.

# State-Value Function $V_\pi(s)$

**Definition:** Discounted return.

- $U_t = R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \cdots + \gamma^{n-t} R_n.$

**Definition:** Action-value function.

- $Q_\pi(s_t, a_t) = \mathbb{E}\left[\, U_t \mid S_t = s_t, A_t = a_t \,\right].$

**Definition:** State-value function.

- $V_\pi(s_t) = \mathbb{E}_A\left[Q_\pi(s_t, A)\right]$

# Action-Value Functions $Q(s, a)$

**Definition:** Discounted return (aka cumulative discounted future reward).

- $U_t = R_t + \gamma\, R_{t+1} + \gamma^2\, R_{t+2} + \gamma^3\, R_{t+3} + \cdots$

**Definition:** Action-value function for policy $\pi$.

- $Q_\pi(s_t, a_t) = \mathbb{E}\,[U_t | S_t = s_t, A_t = a_t].$

**Definition:** Optimal action-value function.

- $Q^\star(s_t, a_t) = \max_\pi Q_\pi(s_t, a_t).$

- Whatever policy function $\pi$ is used, the result of taking $a_t$ at state $s_t$ cannot be better than $Q^\star(s_t, a_t)$.
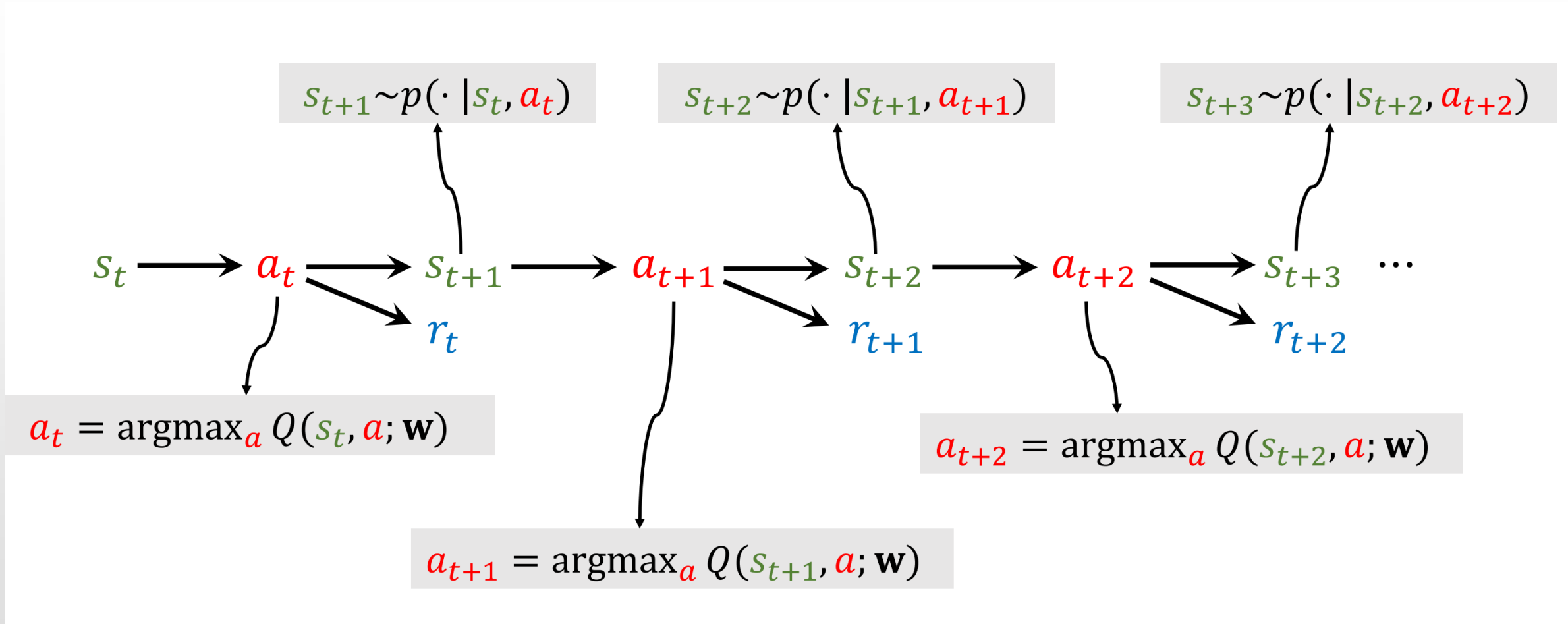
# Approximate the Q Function

**Goal:** Win the game ($\approx$ maximize the total reward.)

**Question:** If we know $Q^\star(s, a)$, what is the best action?

- Obviously, the best action is $a^\star = \underset{a}{\mathrm{argmax}}\, Q^\star(s, a)$.

**Challenge:** We do not know $Q^\star(s, a)$.

- Solution: Deep Q Network (DQN)
- Use neural network $Q(s, a; \mathbf{w})$ to approximate $Q^\star(s, a)$.

# How to apply TD learning to DQN?

**Identity:** $U_t = R_t + \gamma \cdot U_{t+1}$.

**TD learning for DQN:**

- DQN's output, $Q(s_t, a_t; \mathbf{w})$, is an estimate of $U_t$.

- DQN's output, $Q(s_{t+1}, a_{t+1}; \mathbf{w})$, is an estimate of $U_{t+1}$.

- Thus, $\underbrace{Q(s_t, a_t; \mathbf{w})}_{\text{Prediction}} \approx \underbrace{r_t + \gamma \cdot Q(s_{t+1}, a_{t+1}; \mathbf{w})}_{\text{TD target}}$.

# Train DQN using TD learning

- Prediction: $Q(s_t, a_t; \mathbf{w}_t)$.

- TD target:

$$y_t = r_t + \gamma \cdot Q(s_{t+1}, a_{t+1}; \mathbf{w}_t)$$

$$= r_t + \gamma \cdot \max_a Q(s_{t+1}, a; \mathbf{w}_t).$$

- Loss: $L_t = \frac{1}{2}[Q(s_t, a_t; \mathbf{w}) - y_t]^2$.

- Gradient descent: $\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \cdot \frac{\partial L_t}{\partial \mathbf{w}}\Big|_{\mathbf{w}=\mathbf{w}_t}$.
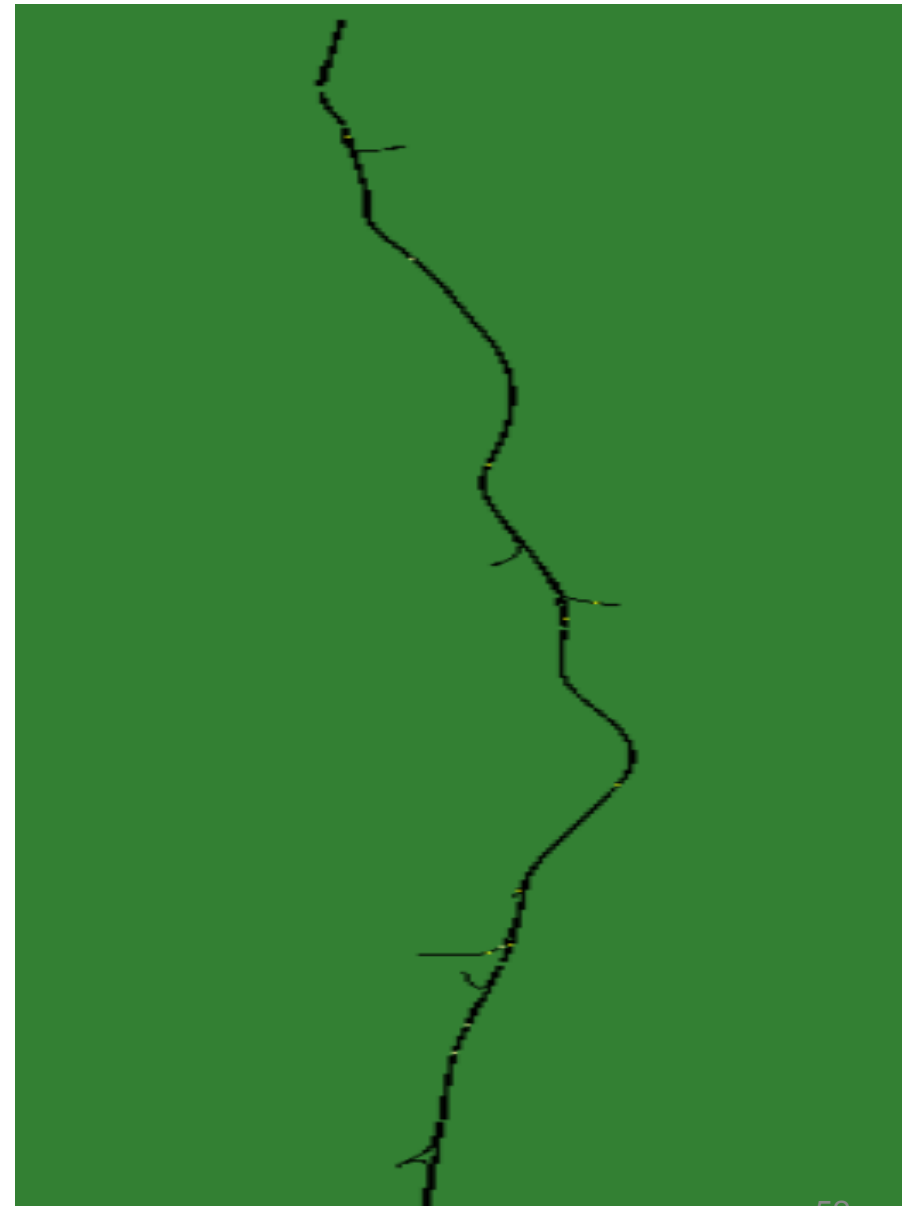
# Temporal Difference (TD) Learning

**Algorithm:** One iteration of TD learning.

1. Observe state $S_t = s_t$ and perform action $A_t = a_t$.

2. Predict the value: $q_t = Q(s_t, a_t; \mathbf{w}_t)$.

3. Differentiate the value network: $\mathbf{d}_t = \dfrac{\partial\, Q(s_t, a_t; \mathbf{w})}{\partial\, \mathbf{w}} \Big|_{\mathbf{w}=\mathbf{w}_t}$.

4. Environment provides new state $s_{t+1}$ and reward $r_t$.

5. Compute TD target: $y_t = r_t + \gamma \cdot \max_a Q(s_{t+1}, a; \mathbf{w}_t)$.

6. Gradient descent: $\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \cdot (q_t - y_t) \cdot \mathbf{d}_t$.
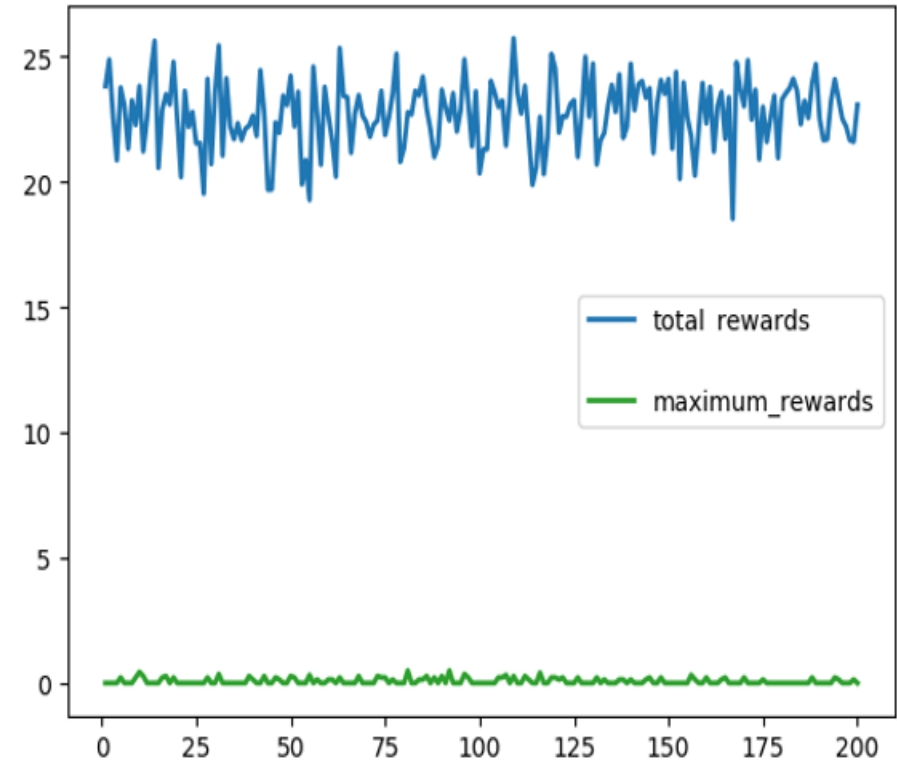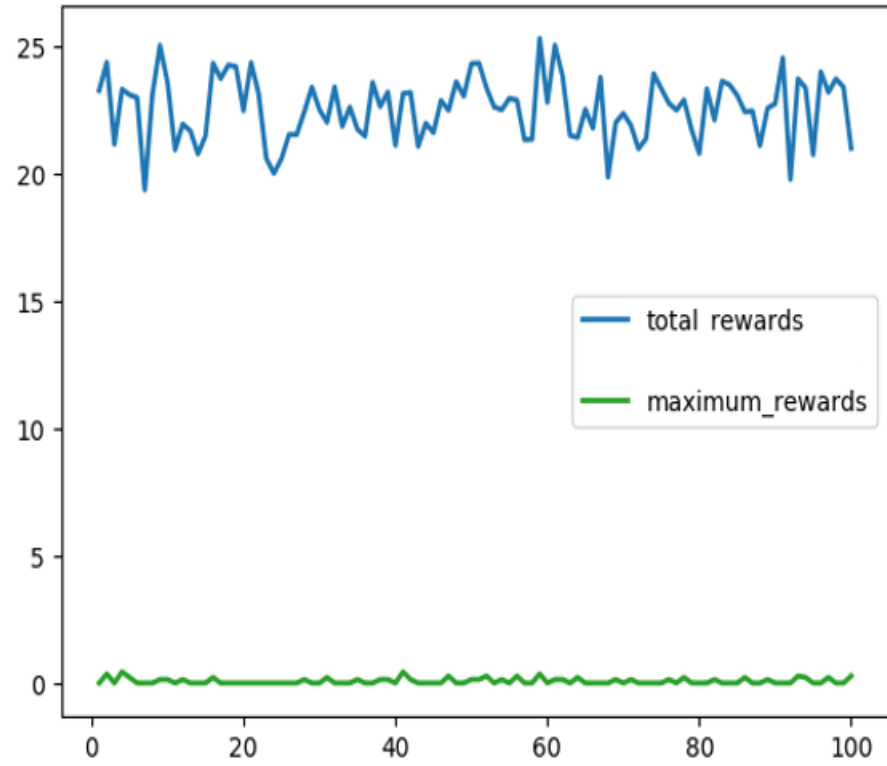
Test: $\min \sum\limits_{t=1}^{T}\sum\limits_{j=1}^{J}\sum\limits_{i=1}^{I} w_{ijt}$

$w_{ijt}$ : the waiting time of vehicle i on ramp j
for time t

$T = 7200, \quad J = 7$

# Result:

**ACKNOWLEDGEMENTS**