

THE CHALLENGE

Strict latency constraints
Edge devices must process continuous sensor data in real-time; cloud offloading introduces unacceptable delays.

Static pipelines fail
Existing systems do not adapt inference mode to workload pressure, causing deadline violations under load.

No unified approach
Frameworks treat latency, energy, and accuracy as separate objectives rather than a joint problem.

RESEARCH GAPS

GAP 1
No explicit latency-aware optimisation for real-time behaviour prediction at the edge.

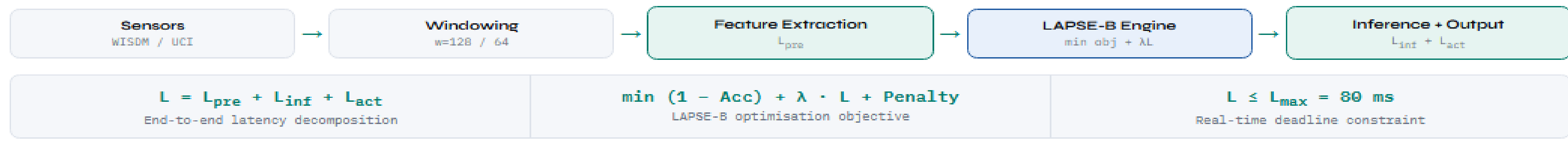
GAP 2
Latency and resource constraints treated separately – no holistic cost function.

GAP 3
Decision overhead of the optimiser itself rarely modelled or bounded.

KEY CONTRIBUTIONS

- Novel latency-aware scheduling algorithm with decomposed end-to-end latency model
- Multi-objective cost function: $(1 - \text{Acc}) + \lambda \cdot L$ with quadratic deadline penalty
- WISDM & UCI HAR dataset integration for realistic HAR workloads
- Real-time demonstration module + interactive analytics dashboard
- Simulation-based evaluation – reproducible, hardware-agnostic

SYSTEM ARCHITECTURE & LAPSE-B DECISION ENGINE



KEY METRICS

- 51.8 ms** LAPSE-B mean latency
- 4.9%** Deadline violations
- 56.4%** Std mode accuracy
- 80 ms** L_{max} deadline

INFERENCE MODES

Standard mode

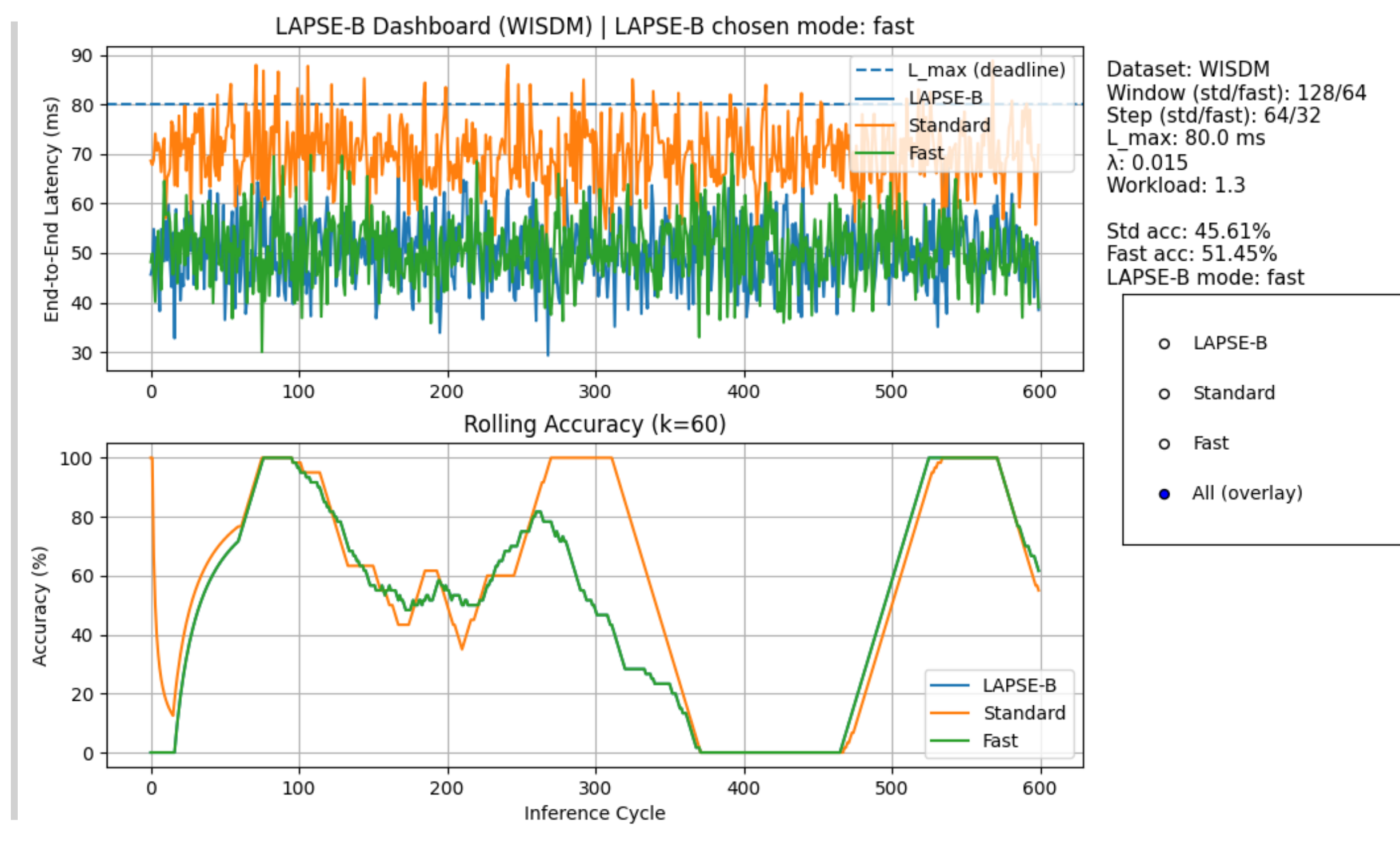
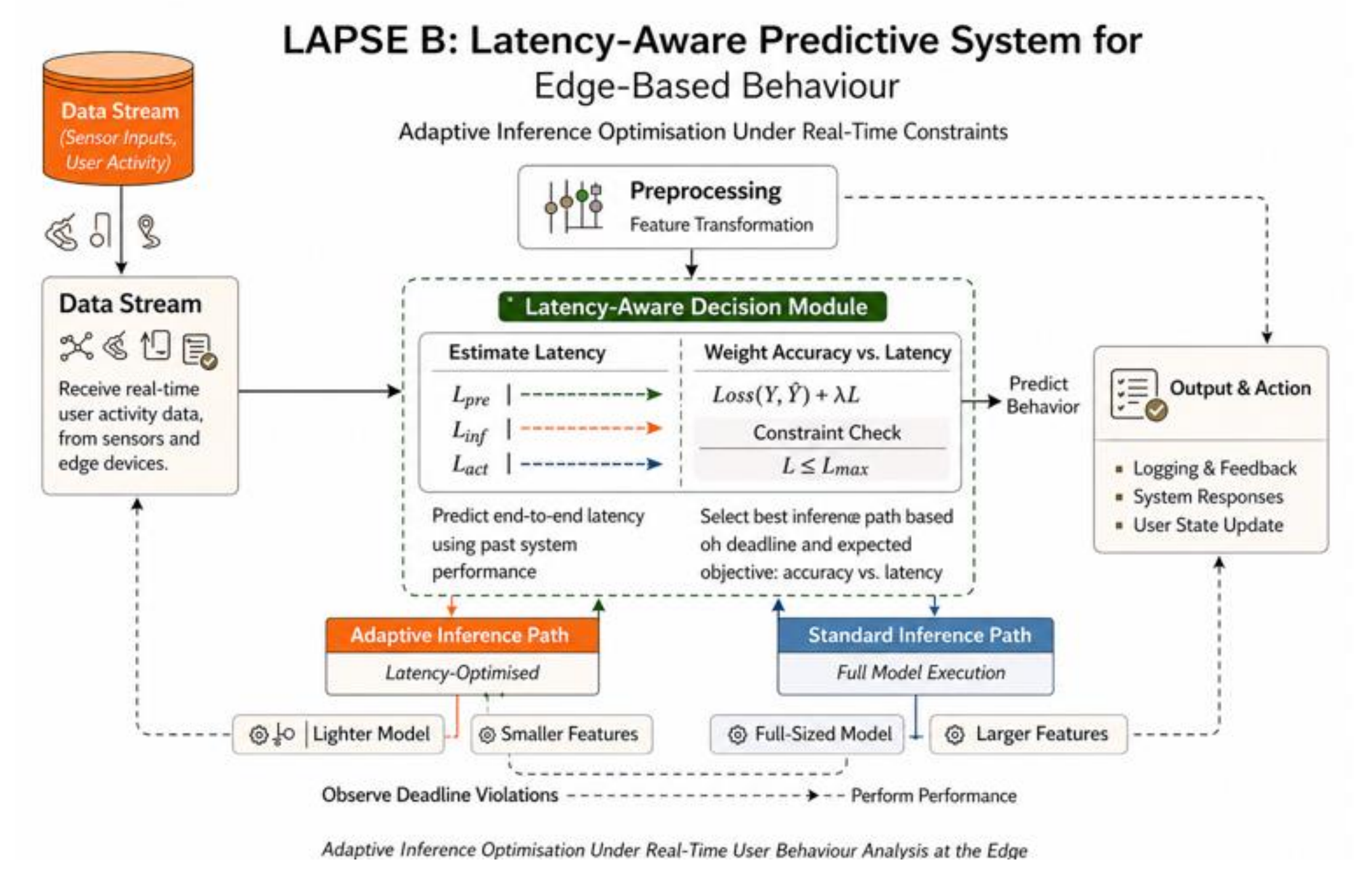
- Window: 128 samples
- Accuracy: 56.43%
- Latency: 71.2 ms
- Violations: 14.3%

Fast mode

- Window: 64 samples
- Accuracy: 49.94%
- Latency: 55.2 ms
- Violations: 6.4%

LAPSE-B (adaptive)

- Mode: Dynamic
- Accuracy: ~49.94%
- Latency: 51.8 ms
- Violations: 4.9%



KEY FINDINGS

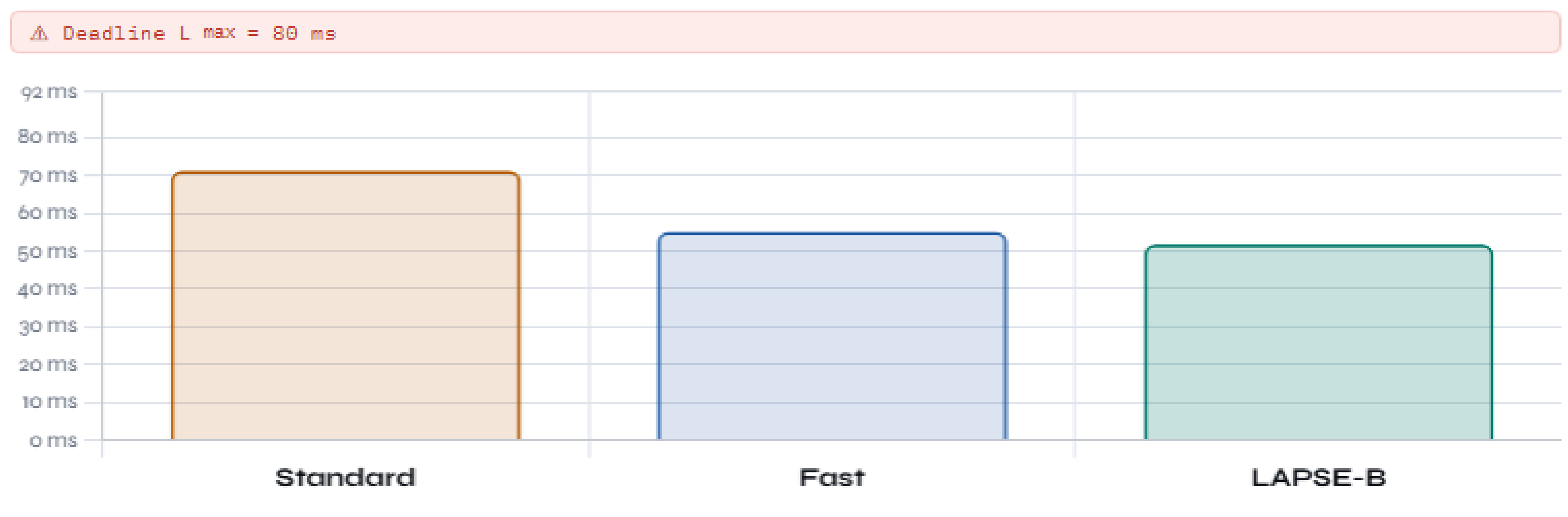
- Latency reduced:** 27% at $\lambda=1.0$, 28% at $\lambda=1.3$, 32% at $\lambda=1.6$ vs Standard
- Violations halved:** 14.3% (Static) → 4.9% (LAPSE-B) under medium load
- Pareto-efficient:** Controlled accuracy degradation – not a naive speed-accuracy swap
- Smooth degradation:** Continuous objective avoids oscillatory mode-switching

CONCLUSION

LAPSE-B demonstrates that **explicit latency awareness** within the inference decision process materially improves real-time performance at the edge – without catastrophic accuracy collapse.

Future work: real hardware deployment on TinyML/mobile platforms; energy-aware objective extension; reinforcement learning for scheduling policy.

LATENCY DISTRIBUTION (MEDIUM WORKLOAD $\lambda=1.3$)



Key result: LAPSE-B reduces deadline violations from 14.3% → 4.9% – a **66% reduction** under medium load ($\lambda=1.3$). Predictive scheduling pre-empts violations rather than reacting to them.

Pareto efficiency: LAPSE-B achieves lower latency than both baselines without proportional accuracy loss – demonstrating controlled, structured degradation rather than a naive speed-accuracy swap.

WORKLOAD ESCALATION

λ	LAPSE-B	Std	Fast
1.0	42.1 ms	65.3 ms	50.4 ms
1.3	51.8 ms	71.2 ms	55.2 ms
1.6	59.7 ms	88.4 ms	63.5 ms

SAVING vs STD ($\lambda=1.6$)
LAPSE-B saves **28.7 ms**

Std exceeds deadline **+8.4 ms**

Legend: WISDM, UCI HAR