

# Small Language Model (SLM) for Device AI

**Akraino Robotics Blueprint, Release 8 Enhancement**



# Quick Background

## › Akraino Robotics Blueprint

- › Led by Fujitsu and Univ Ritsumeikan, incubation 2022
- › Sponsored by SIP/Japan Cabinet Office / NEDO
- › Signalogic added real-time ASR <sup>1</sup>

## New industries for robots



Agriculture



Food factory



Retail



Restaurant

## Challenges for robot in these industries

1. Objects with diverse shapes, flexibility
2. Uncertain environment (wet, clutter, customers, etc.)
3. “No cloud” communication with humans

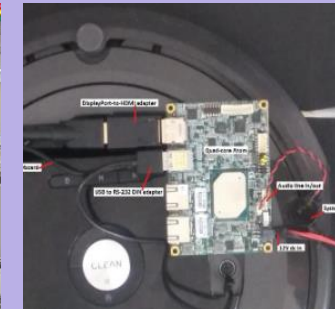


### *Solutions through fusion of robot and sensors*

- Flexible robot handling using sensor data
- Reliable, light weight, onboard ASR <sup>1</sup>

“CPS <sup>2</sup> Robot blueprint family” published as OSS stack in Akraino

PoC in progress  
in the field



<sup>1</sup> Automatic Speech Recognition

<sup>2</sup> CyberPhysical Systems

# Device AI speech recognition challenges at the edge

- › Device AI applications need to run ASR
  - › With unreliable or no cloud connection
  - › On very small form-factor devices (e.g. pico ITX)
  - › Under difficult conditions, including background noise, urgent or stressed voice input, and background talkers
  - › Robotics servo motor and other mechanical noise increases difficulty



# Precise Command Problem

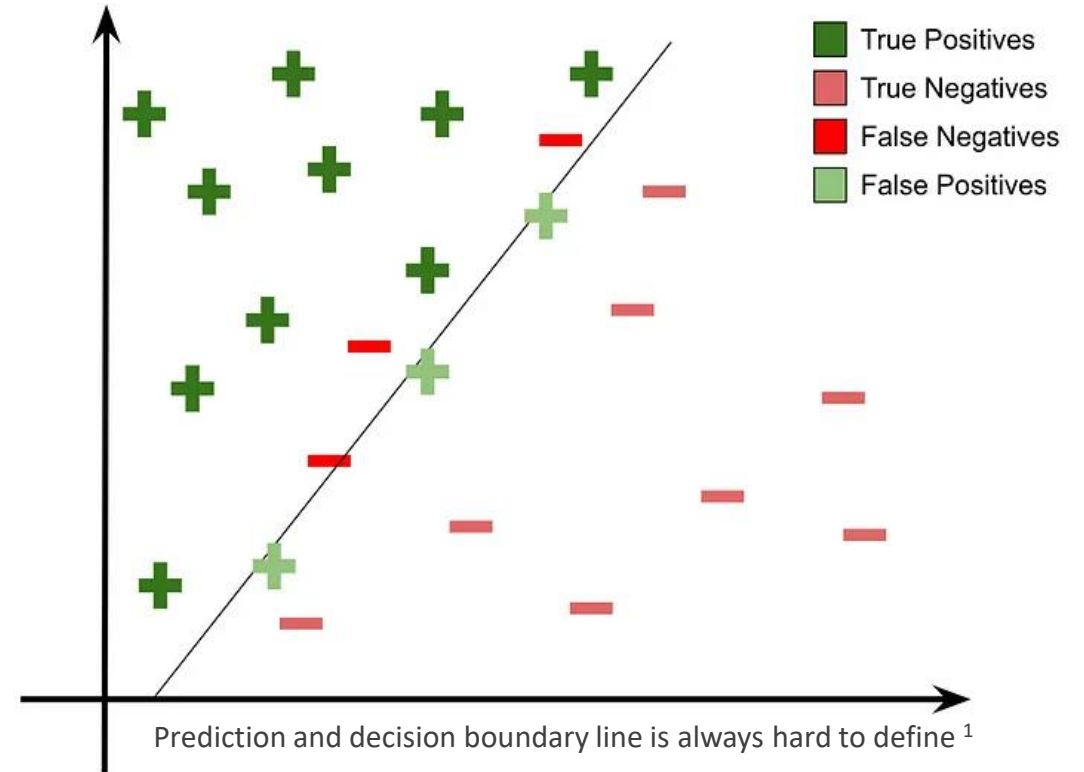
- › Machine-readable APIs must be precise
- › Minimizing false positives is crucial
- › Under difficult conditions, efficient open source ASRs such as Kaldi and Whisper produce "sound-alike" errors, for example:

`"in the early days a king rolled the stake"`

which must be corrected to

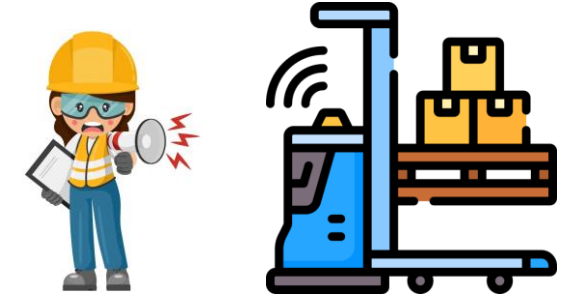
`"in the early days a king ruled the state"`

- › Sound-alike errors are problematic for safety and emergency situations
  - › Internet / cloud connection cannot be assumed. Phones may be useless
  - › A first responder may use a portable hand-held device and give commands to a robotaxi such as "get off the road in that turn-out up ahead and shut down"
  - › Sometimes generalized in ASR research as "substitution errors"



# Use Cases

- › Factory floor personnel need to give urgent commands
  - › possibly dangerous equipment (e.g. forklifts)
  - › no-hands-free environments (e.g. food processing)
- › First responders need to communicate with disabled or disconnected robotic vehicles
  - › robotaxis
  - › semi trucks
- › Language Translation
  - › sound-alike correction prior to translation
  - › independent of ASR model



# Requirements

- › Must correct sound-alike errors independently of ASR model without re-training, tuning, compression, or other reduction
- › Very small form-factor, under 15 W
  - › 4 x 3", heat sink only, no fans
- › Real-time – must run every 250 to 500 msec
  - › Minimum 10 token/sec, preferably 20
- › Backwards / forwards context of 5 tokens (words)
  - › unlike an LLM, wide context window, domain knowledge, and extensive web page training are not needed
- › Compliant with emerging teleoperation standards
  - › California included teleoperation as part of its regulation for driverless vehicles in 2018
  - › NIST conference in 2020
  - › WiFi or USB port interfaces typical

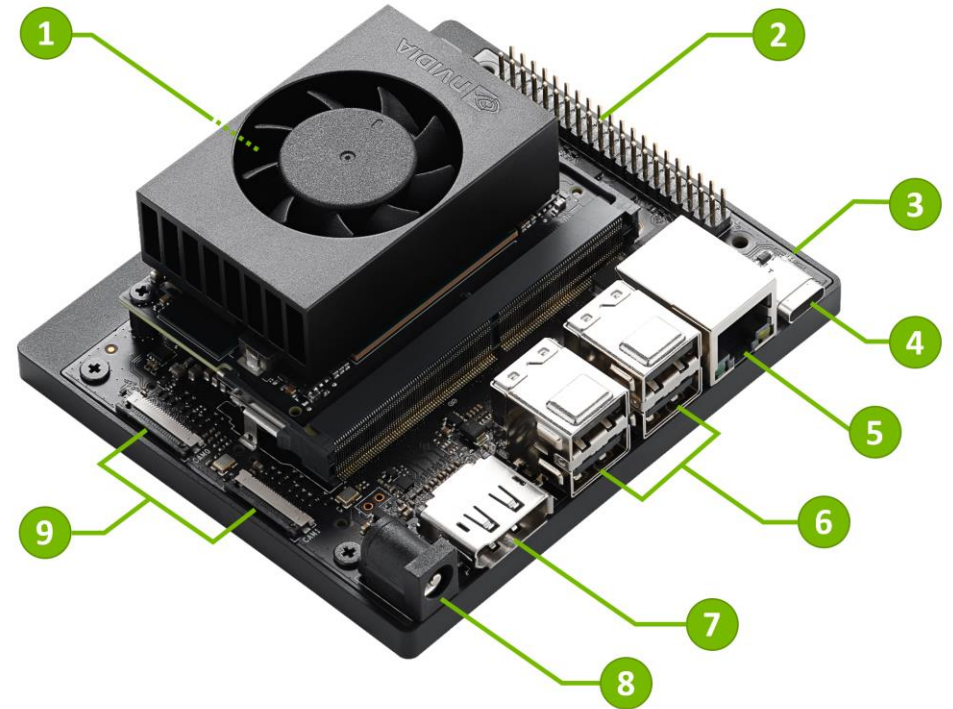


Teleoperation and Autonomous Vehicles Overview		
	Key Information	Other Information
What is teleoperation?	<ul style="list-style-type: none"> <li>• Remote operation of a machine at a distance</li> <li>• Requires wireless link to machine</li> <li>• First concepts in 1870s: wire-guided torpedoes</li> </ul>	<ul style="list-style-type: none"> <li>• Similar to remote control</li> <li>• Or wired link if machine is nearby</li> <li>• Nikola Tesla-1898: Radio-controlled boat</li> </ul>
3 levels of AV teleoperation	<ul style="list-style-type: none"> <li>• Remote monitoring of AVs</li> <li>• Remote assistance to AVs</li> <li>• Remote driving of AVs</li> </ul>	<ul style="list-style-type: none"> <li>• Monitoring of AV fleet driving</li> <li>• Driving assist for a short time</li> <li>• Driving for a substantial time</li> </ul>
Why is it needed?	<ul style="list-style-type: none"> <li>• As human backup to driverless vehicles</li> <li>• To manage and learn from edge cases</li> <li>• To gain early AV deploy with acceptable safety</li> </ul>	<ul style="list-style-type: none"> <li>• To be part of most AV regulations</li> <li>• Transfer edge cases to known cases</li> <li>• Only for specific AV use-cases</li> </ul>
Teleoperation regulation status	<ul style="list-style-type: none"> <li>• California approval granted in February 2018</li> <li>• California operational use started in April 2018</li> <li>• Countries: Canada, Finland, Japan, Netherlands</li> <li>• Shanghai and other Chinese cities</li> </ul>	<ul style="list-style-type: none"> <li>• Driverless AVs require teleoperation</li> <li>• AZ, FL, MI, OH, TX too; More will follow</li> <li>• Sweden, UK; More will follow</li> <li>• Teleoperation expected in China</li> </ul>
Teleoperation use-cases	<ul style="list-style-type: none"> <li>• Sidewalk AVs: Most common usage</li> <li>• Trucks: AV on highway; last mile teleoperation</li> <li>• Robotaxis: Regulation and edge case</li> <li>• Others: Forklifts, excavators, yard trucks, combine</li> <li>• Shared electric scooters</li> </ul>	<ul style="list-style-type: none"> <li>• Examples: Kiwibot, Postmates</li> <li>• Examples: Einride, Hub-to-hub AVs</li> <li>• Zoox has remote operation patent</li> <li>• Testing, trials, some deployment</li> <li>• To return to base &amp; charging stations</li> </ul>
Teleoperation startups	<ul style="list-style-type: none"> <li>• Designated Driver: Assisted &amp; remote driving</li> <li>• DriveU: Assisted &amp; remote driving teleoperation</li> <li>• Ottopia: Assisted &amp; remote-driving teleoperation</li> <li>• Phantom Auto: Focus on remote driving use-cases</li> </ul>	<ul style="list-style-type: none"> <li>• Teleoperation for Texas A&amp;M shuttle</li> <li>• Member: Israeli teleoperation consortium</li> <li>• Partners: BMW, Denso, EasyMile, others</li> <li>• Forklifts, yard trucks and similar clients</li> </ul>
Make or buy teleoperation?	<ul style="list-style-type: none"> <li>• Top AV software platform: own teleoperation</li> <li>• Many companies will buy teleoperation software</li> </ul>	<ul style="list-style-type: none"> <li>• Likely integrated with AV software driver</li> <li>• From multiple teleoperation startups</li> </ul>
Teleoperation standards	<ul style="list-style-type: none"> <li>• Teleoperation standards likely to happen</li> <li>• Best chance is high level standards</li> </ul>	<ul style="list-style-type: none"> <li>• AV software driver variety is big barrier</li> <li>• At functional or operational level</li> </ul>
Teleoperation Forum	<ul style="list-style-type: none"> <li>• First conference on teleoperation (virtual)</li> <li>• NIST Vehicle Teleoperation Forum   NIST</li> </ul>	<ul style="list-style-type: none"> <li>• November 13, 2020 by NIST</li> <li>• 40 speakers; 8+ hours of video sessions</li> </ul>
Teleoperation Consortium	<ul style="list-style-type: none"> <li>• TC is a non-profit business organization</li> <li>• 30+ companies, universities, organizations</li> </ul>	<ul style="list-style-type: none"> <li>• Founded December 2020</li> <li>• Website: Teleoperation Consortium</li> </ul>
NIST=National Institute of Standards and Technology		
Source: Egil Juliusen, May 2021		



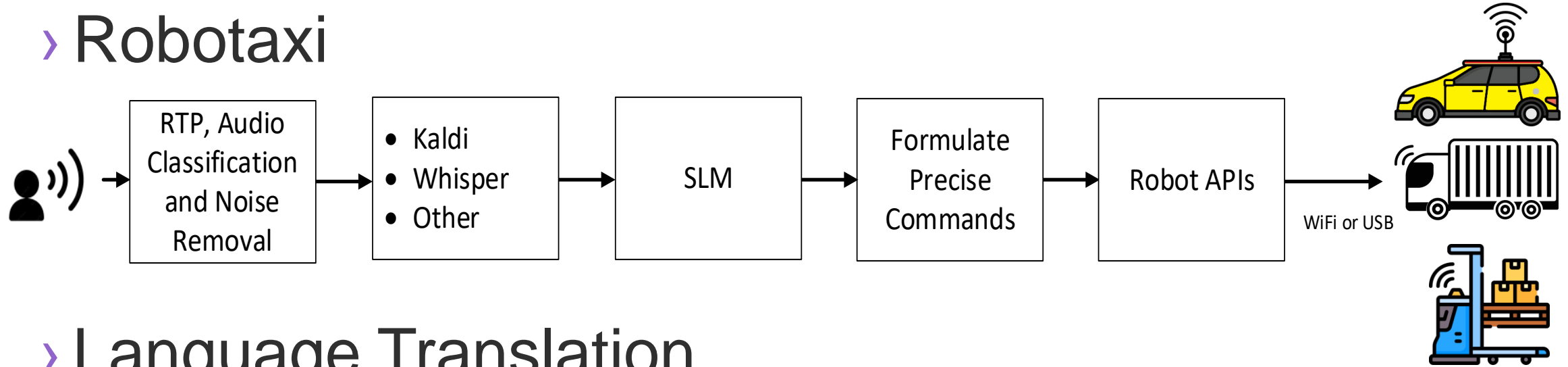
# Why Not Use a GPU ?

- › An Nvidia Jetson Orin Nano running 7B weight Llama-2 model meets power consumption and size requirements but ...
  - › testing shows we can't reliably fix sound alike errors – context required too wide and model too small
  - › processing rate is approx 1 token / sec
  - › we still need to run RTP, audio classification, noise removal, and ASR. These are algorithms and sequential logic, extremely difficult to code in Cuda
  - › we may need additional CPU cores for application-specific requirements
- › Needs a fan, can't operate heat-sink only

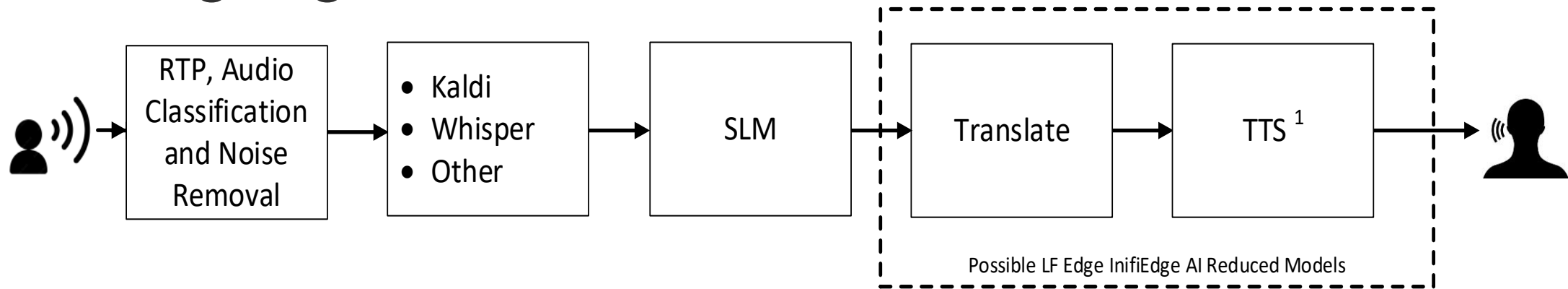


# Technology Overview – Dataflow

## › Robotaxi



## › Language Translation





# Technology Overview – Training and Inference

- › Conventional CPUs
  - › Arm, x86
  - › no GPUs, no HBM
- › Conventional memory, 8 GB min
- › Training
  - › frequency domain representations of 10,000 text words – becomes an image recognition problem
    - › non-linear memory space, self-organizing, sound-alikes are near each other
    - › extremely fast
    - › no gradient descent or other high complexity algorithms
- › Inference
  - › content addressable memory – series of spans and local searches



# Status and Next Steps

- › Working now
  - › RTP, audio classification, noise removal - one Atom core in real-time
  - › Kaldi ASR - one Atom core in real-time
  - › pico ITX board (quad core Atom x5-E3940)
  - › 20,000 word vocabulary
- › SLM under development
  - › live demo next step
  - › pico ITX board
  - › planning for Akraino Fall Summit

