



System Engineering Aspects of End-to-End tinyML

Jan Ernst,
Director of AI

Latent AI
November 2021

Latent AI – Offerings

Software to optimize AI models for low power and latency, faster time to market.

Deployment workflows for various hardware.

State of the art reference models for common tasks.

Tooling throughout ML lifecycle: from design exploration to enterprise deployment management.

White paper: <https://latentai.com/unlocking-the-power-of-edge-ai/>

Awards



<https://bit.ly/2EVXsBF>



Top 100 AI Startup 2021



<https://tcn.ch/34juQee>

<https://vimeo.com/460028482>



Top 60 Edge Computing Companies

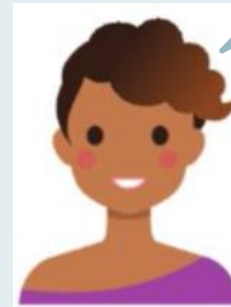
Investors



ML Lifecycle

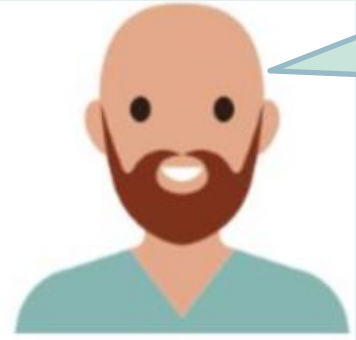
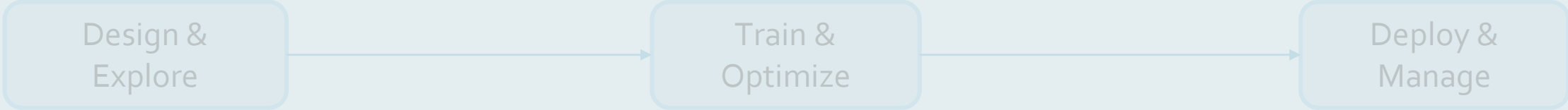


ML Expert Evelyn



I want to build cool stuff!

The Journey of ML Expert Evelyn



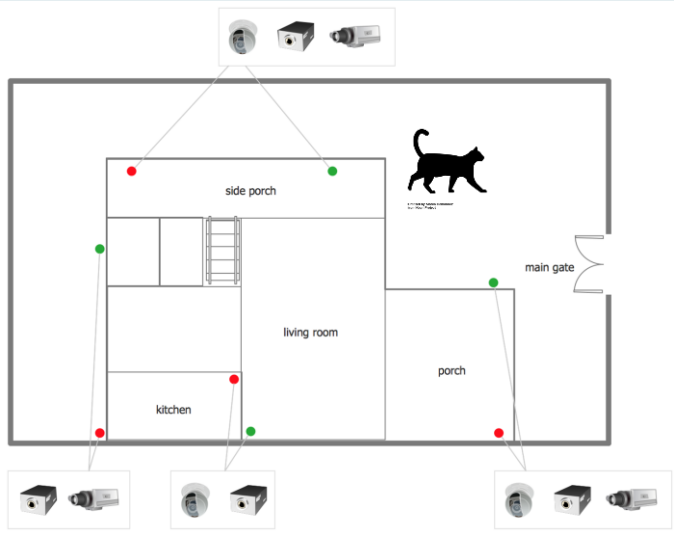
Solutions Steve
(Links with customer)

Let's build a new application:
"Find my cat!"

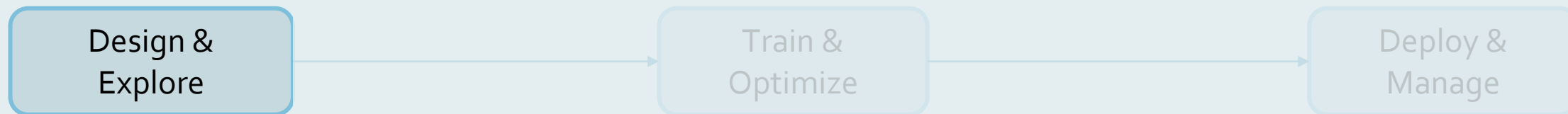
Awesome! Tell me more about the requirements...



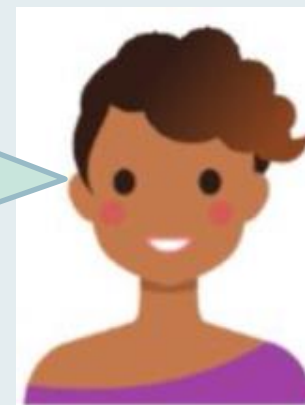
Expert Evelyn



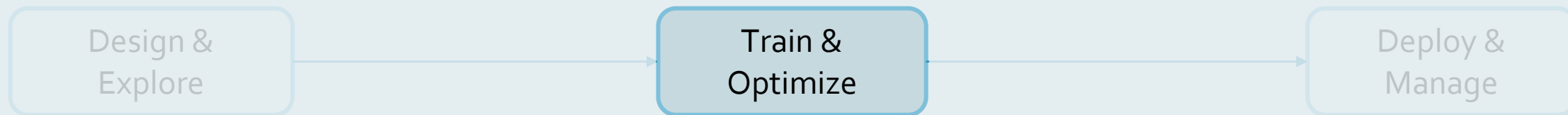
The Journey of ML Expert Evelyn



What hardware fits into my power envelope?
What model fits into my hardware?
What data do I have/need?

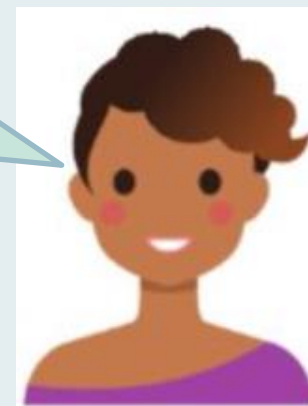


The Journey of ML Expert Evelyn

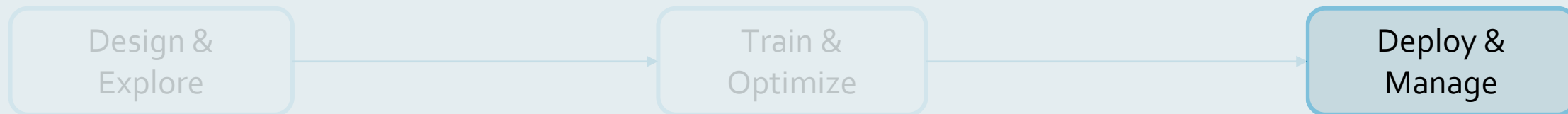


How do I train my model on my data?

How do I optimize my model for the hardware?



The Journey of ML Expert Evelyn



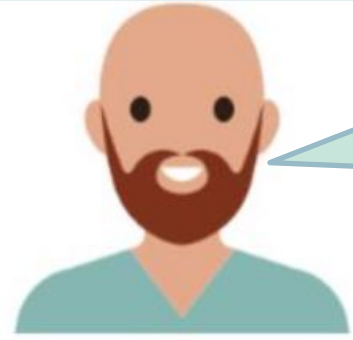
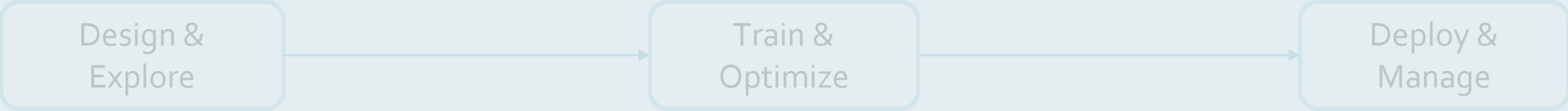
How do I run my model on device?

How do I integrate into a larger system?

What about enterprise reporting?



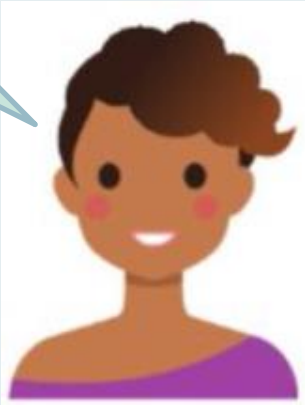
The end of the Journey?



Solutions Steve

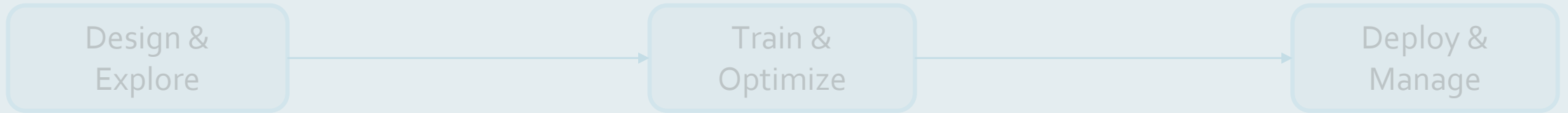
Awesome! I'll start integrating and selling this!

Here is the deployable component!



Expert Evelyn

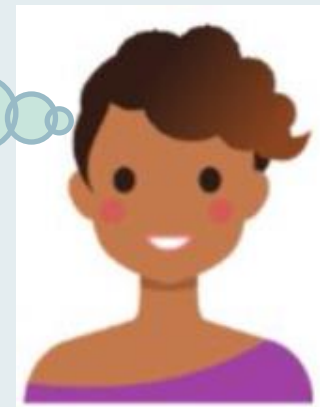
Somewhere down the road



Solutions Steve

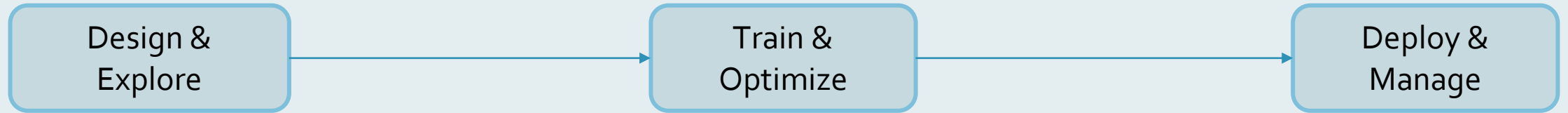
“Find my cat!” is a hit!

What about **dogs**?
Can we add **RFID**?
Need to recognize pet **identity**!
Can I use **cheaper** HW?
Let’s **voice-activate** our app!



Expert Evelyn

Retrospective



Early opportunistic design choices instead of **holistic optimization**

- “Choose model due to framework support”
- “Found nice git repo with model and training”

Evelyn needs best of breed, latest available state of the art from various ecosystems.

Evelyn must develop MxN combinations of model and optimization methods.

Evelyn needs to validate performance criteria of task vs. target HW.

Evelyn must replicate her accuracy evaluation code with the vendor API and guarantee that it's the same performance as in training phase.



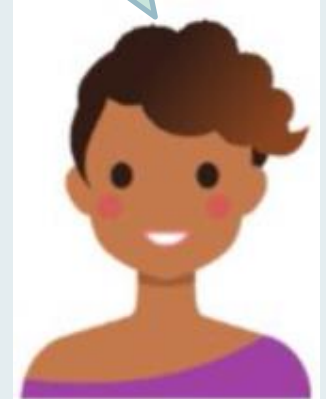
How to improve?

Apply Software Engineering design principles to ML Engineering

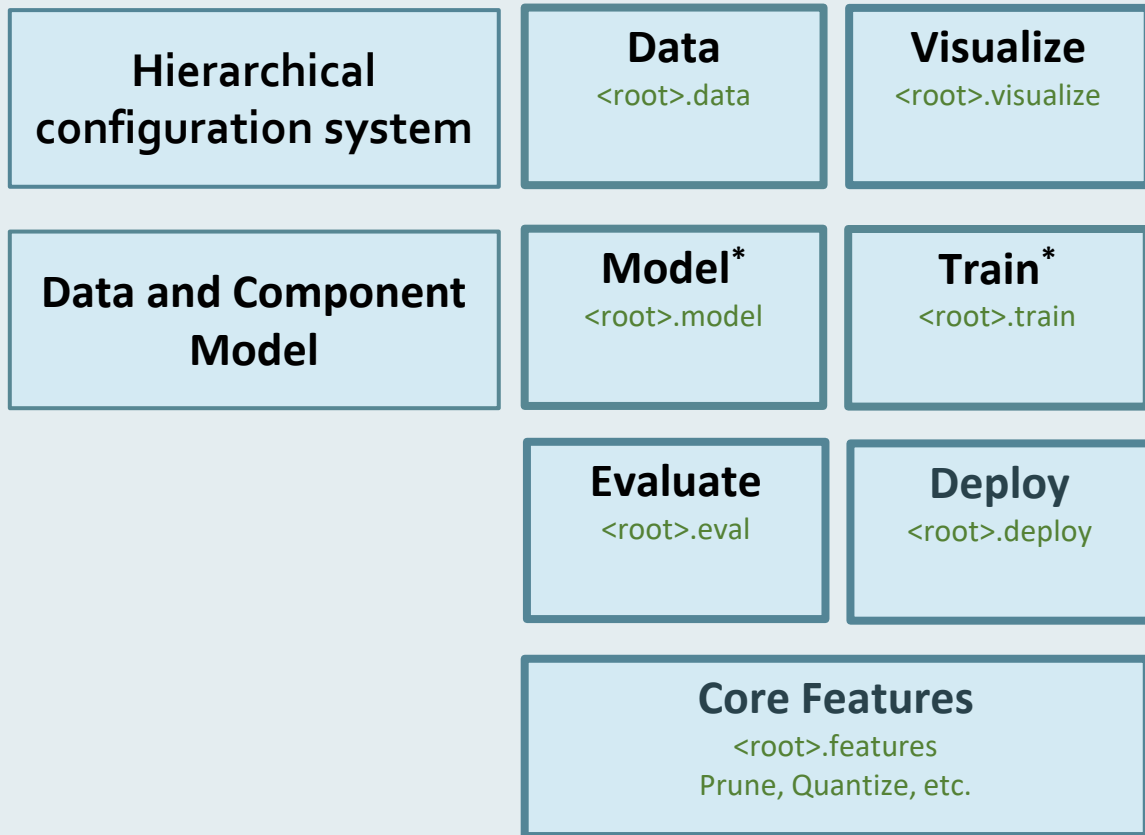
- Compositionality and Interoperability
- Separation of concerns
- Autonomous components

“**System**”: thin, flexible shell of loose coupling of components

What can I learn from the 1960's?



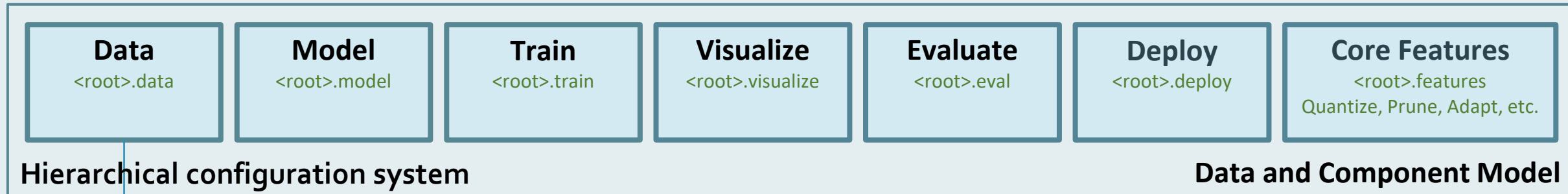
Our approach – Compositional and Modular



A Component is

- Self-contained, independent
- Describes itself completely via interfaces and implied semantics
- Exposes its possible interactions with other components

System Perspective



Expose two types of Interfaces:

“What am I?”

- Images, Signals, Tensors
- Metadata, Annotations

“How do I interact with other components?”

- Augmentation
- Transparent caching

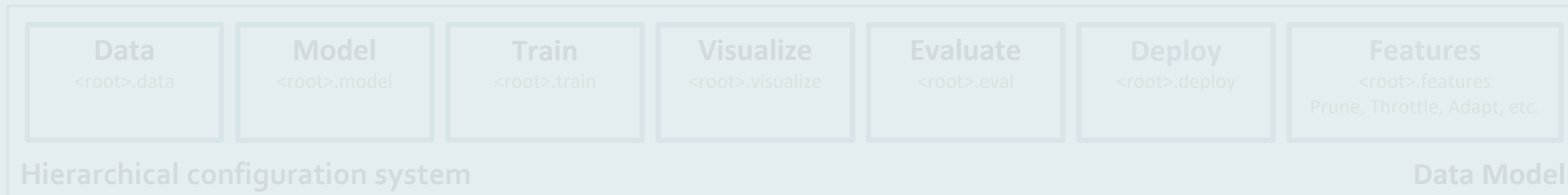
“What does data represent?”

- “Annotation is something a sample may have”
- “Annotations can have confidence, geometry, timestamp, ...”
- “Bounding box is quadrilateral spatial geometry”

“How can data be interpreted?”

- “Object detections can be image labels”
- “Bounding box can be weak supervision for segmentation”

System Perspective



Resulting system is **modular** and **compositional**

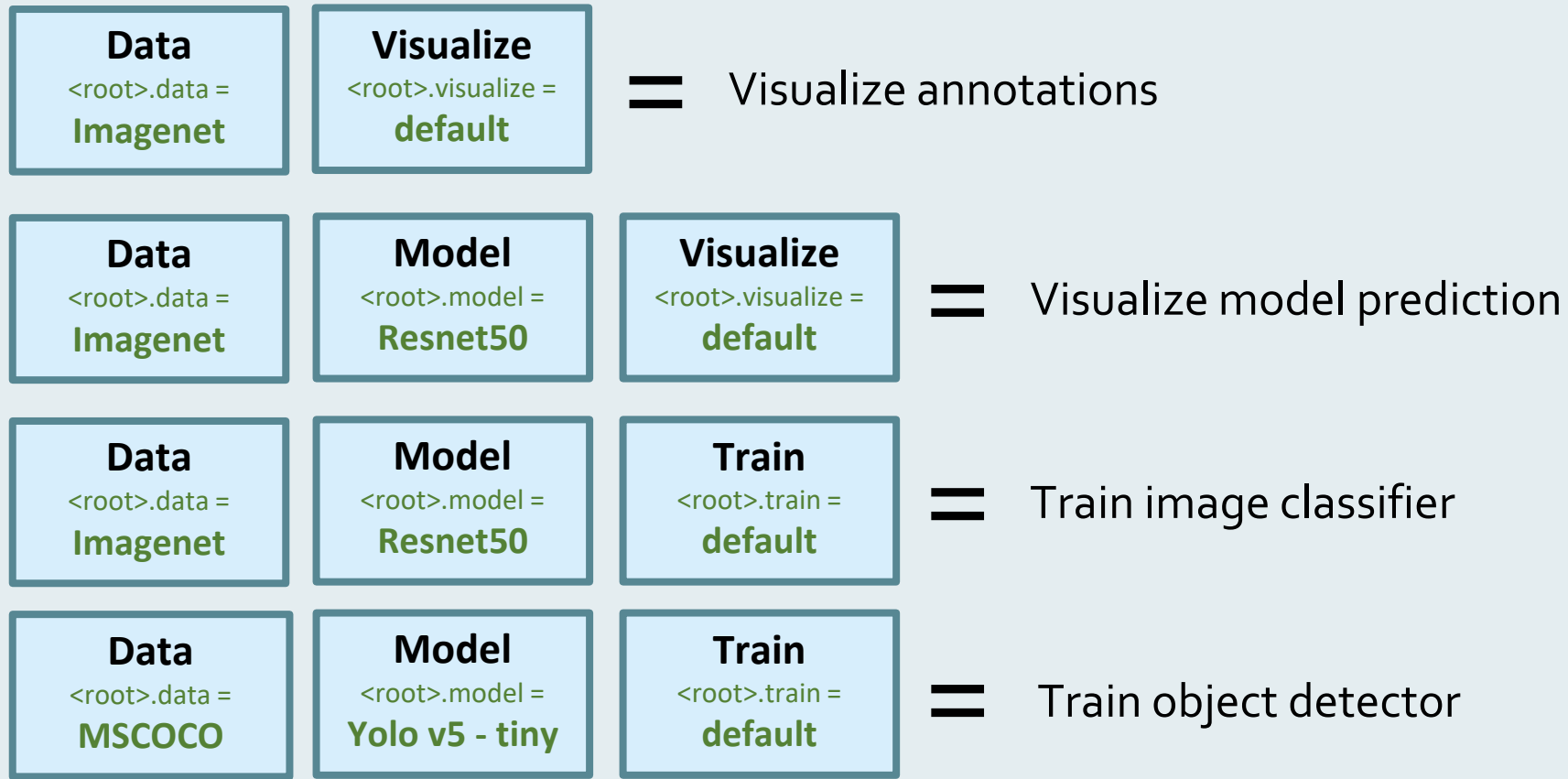
Any component can be developed on its own.

A system is fully described by the set of its components.

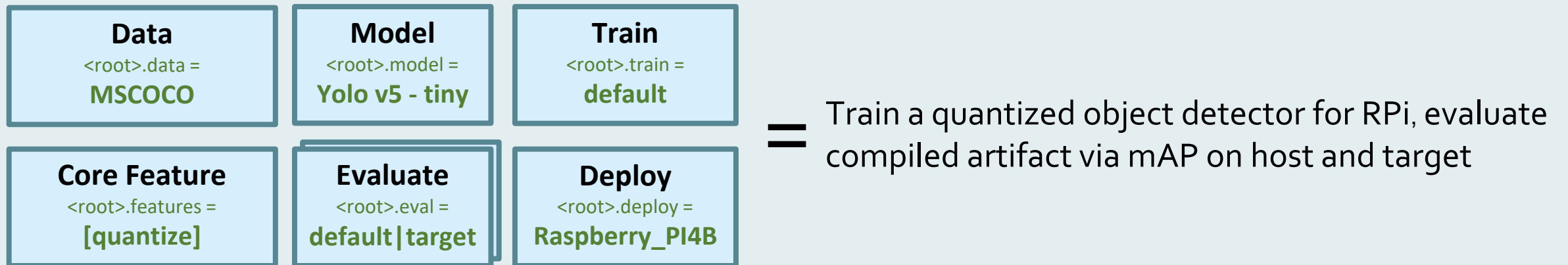
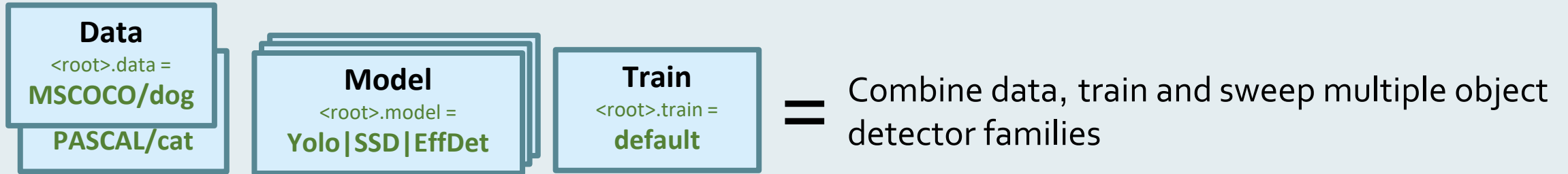
Abstractions not in the way of **going deep anywhere**.

The power of compositionality: No-code system design (only YAML/CLI configuration)

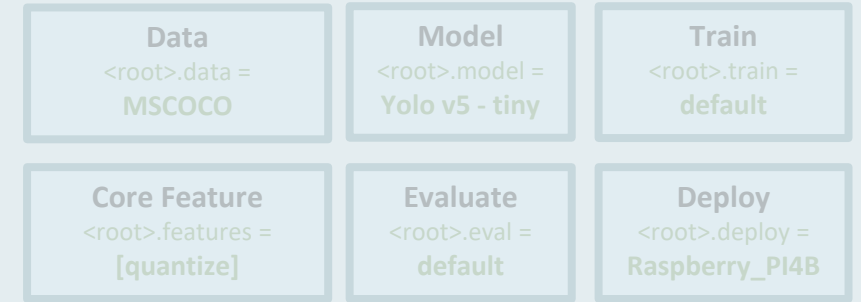
Examples



The power of compositionality: No-code system design (only YAML/CLI configuration)



How does it help Evelyn to scale?



Faster time to market: Can quickly design and sweep new variations and evaluate end to end (data to device).

Simplified ML CICD: Individual components or whole system, don't need expert on end-to-end pipeline.

Simplifies ML Dev: Improve parts without needing to understand whole pipeline: Separate Data, Model, Deploy teams.

Simplifies ML Ops: Reproducible processes, end-to-end tracking of data and model provenance.

Trusted results: Guaranteed same evaluation on host and target HW (preprocessing, postprocessing, protocol)

Best practices built-in or interoperable: distributed, mixed precision, fault tolerance, experiment tracking, etc.

Thank You

Contact us at
info@latent.ai

www.latent.ai

Please fill out a tinyML survey



<https://bit.ly/3GJ5etE>

Latent AI is conducting a short survey to better understand how engineers and developers make design choices for their tiny ML systems. We will share survey results with those that participate.