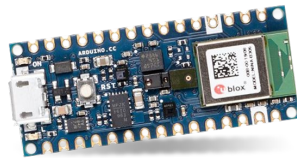
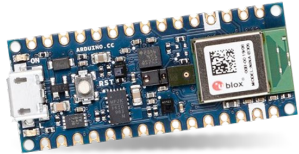
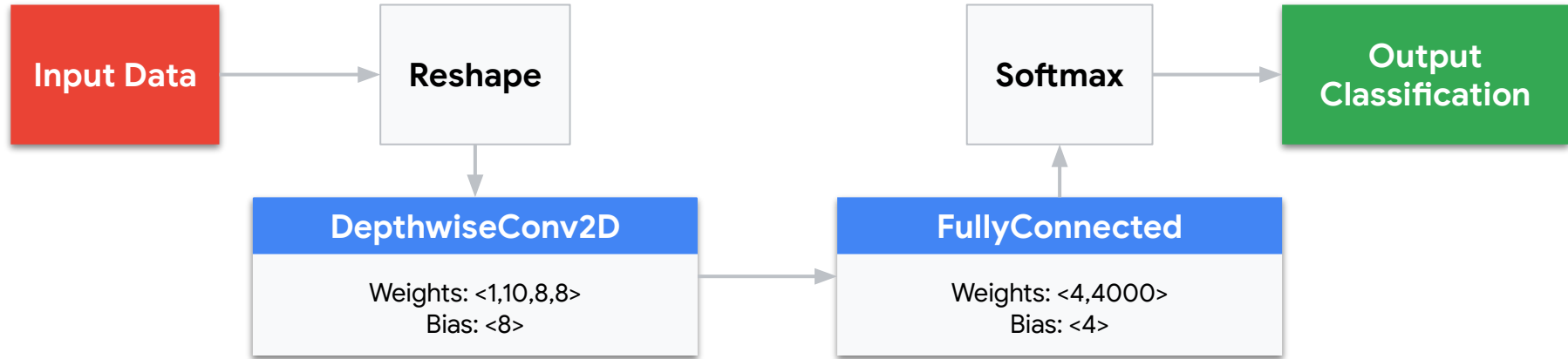


Cascade Architectures

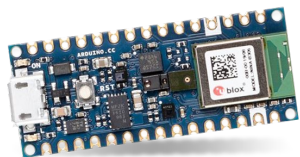
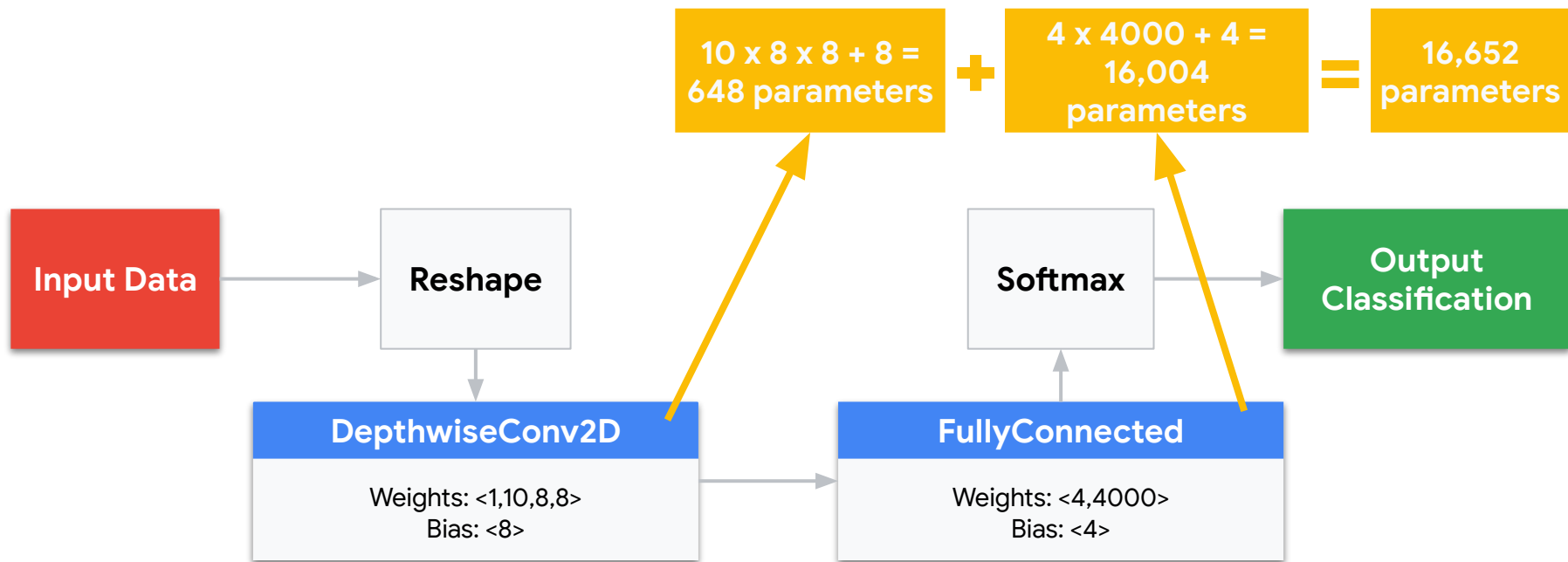




Our board [Course 3 Kit] only has **256KB** of RAM (memory)



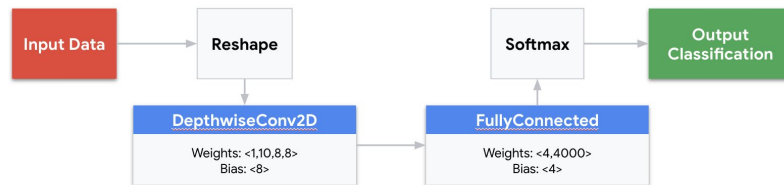
Our board [Course 3 Kit] only has **256KB** of RAM (memory)



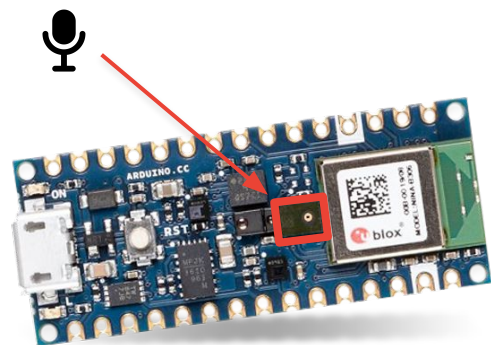
Our board [Course 3 Kit] only has **256KB** of RAM (memory)

Trade-offs

- Limited **vocabulary**
- Limited **accuracy**
- Limited **user experience**



“Cascade” detection: a multi-stage model

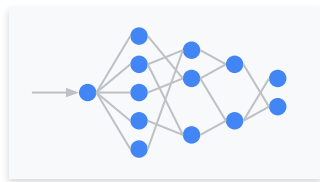


1

Continuously listen on the microcontroller

2

Process the data with **TinyML** at the edge



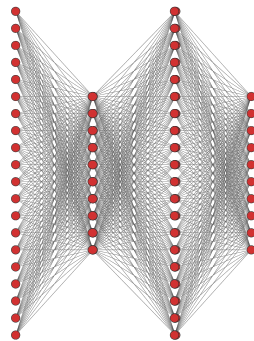
3

Send the data to the cloud when triggered

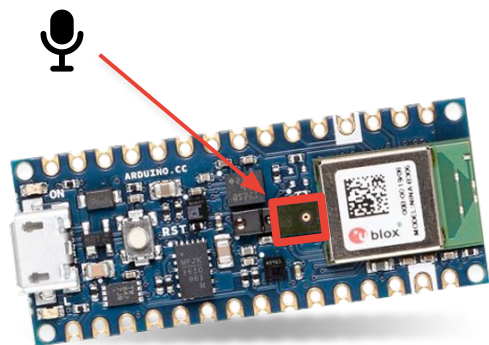


4

Process the full speech data with a large model in the cloud



“Cascade” detection: a multi-stage model

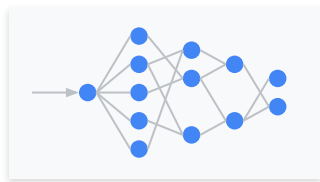


1

Continuously listen on the microcontroller

2

Process the data with **TinyML** at the edge



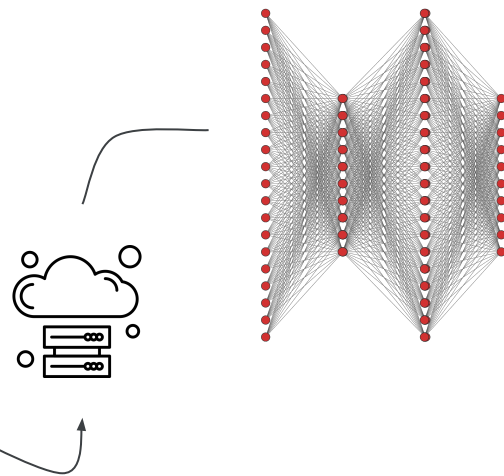
3

Send the data to the cloud when triggered

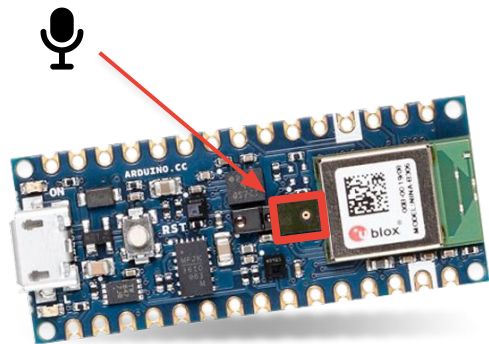


4

Process the full speech data with a large model in the cloud



“Cascade” detection: a multi-stage model

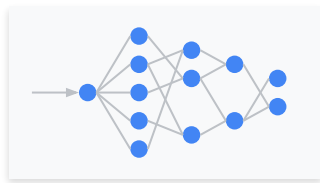


1

Continuously listen on the microcontroller

2

Process the data with **TinyML** at the edge

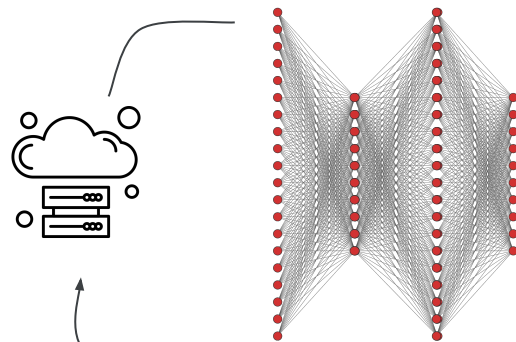


3

Process on a secondary larger model on a larger local device

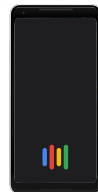
5

Process the full speech data with a large model in the cloud



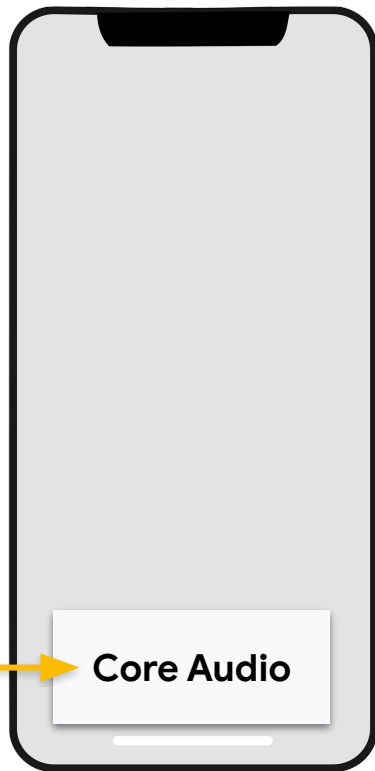
4

Send the data to the cloud when triggered



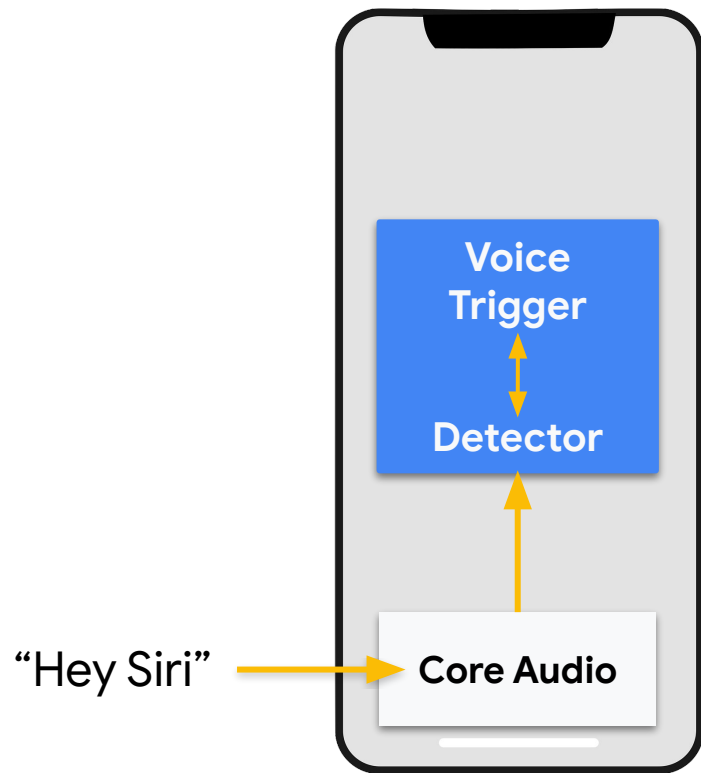
“Hey Siri”

“Hey Siri”

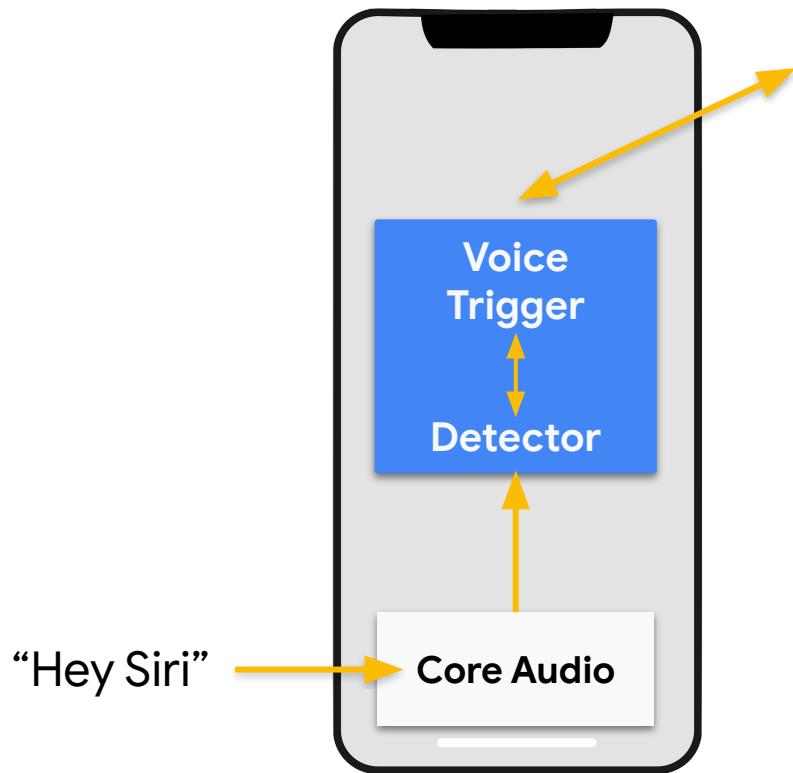


Core Audio

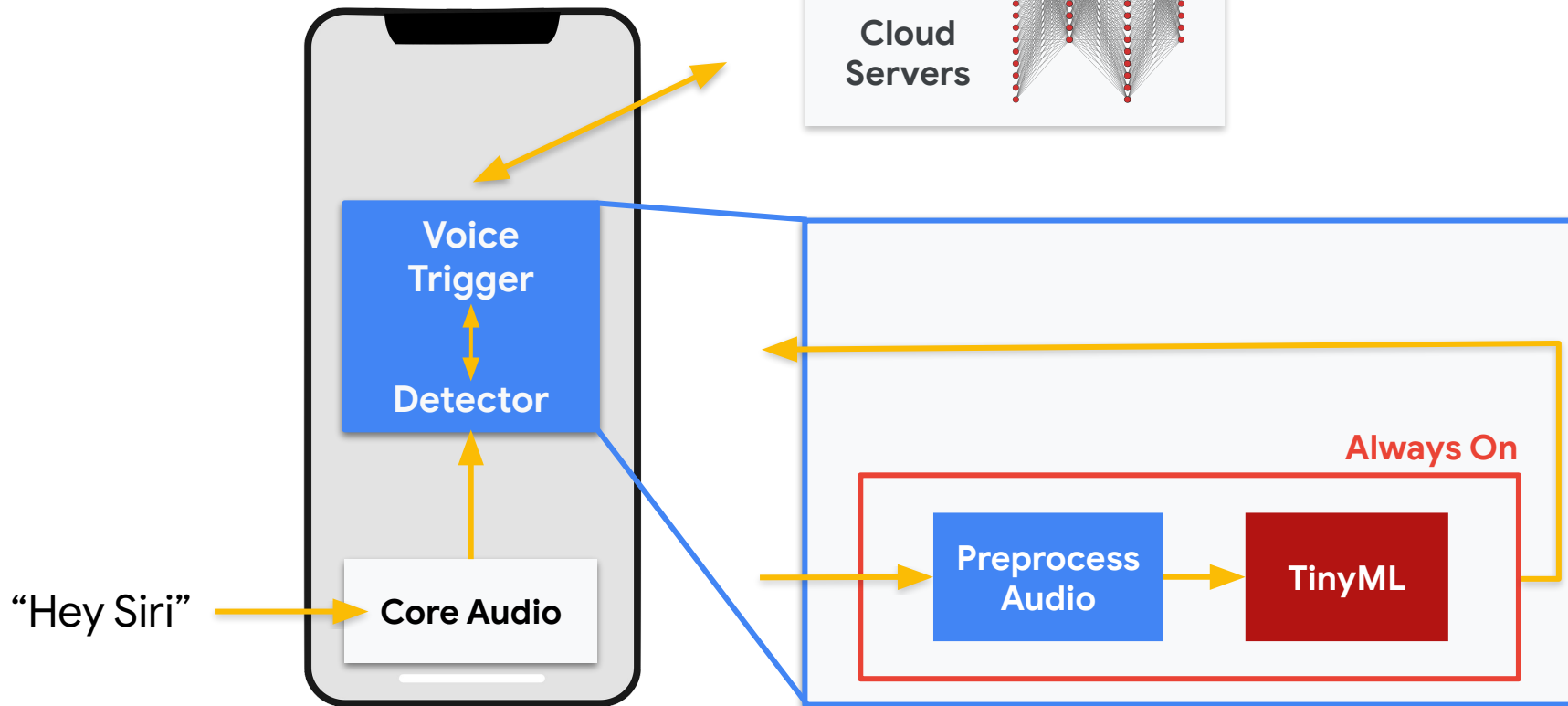
“Hey Siri”



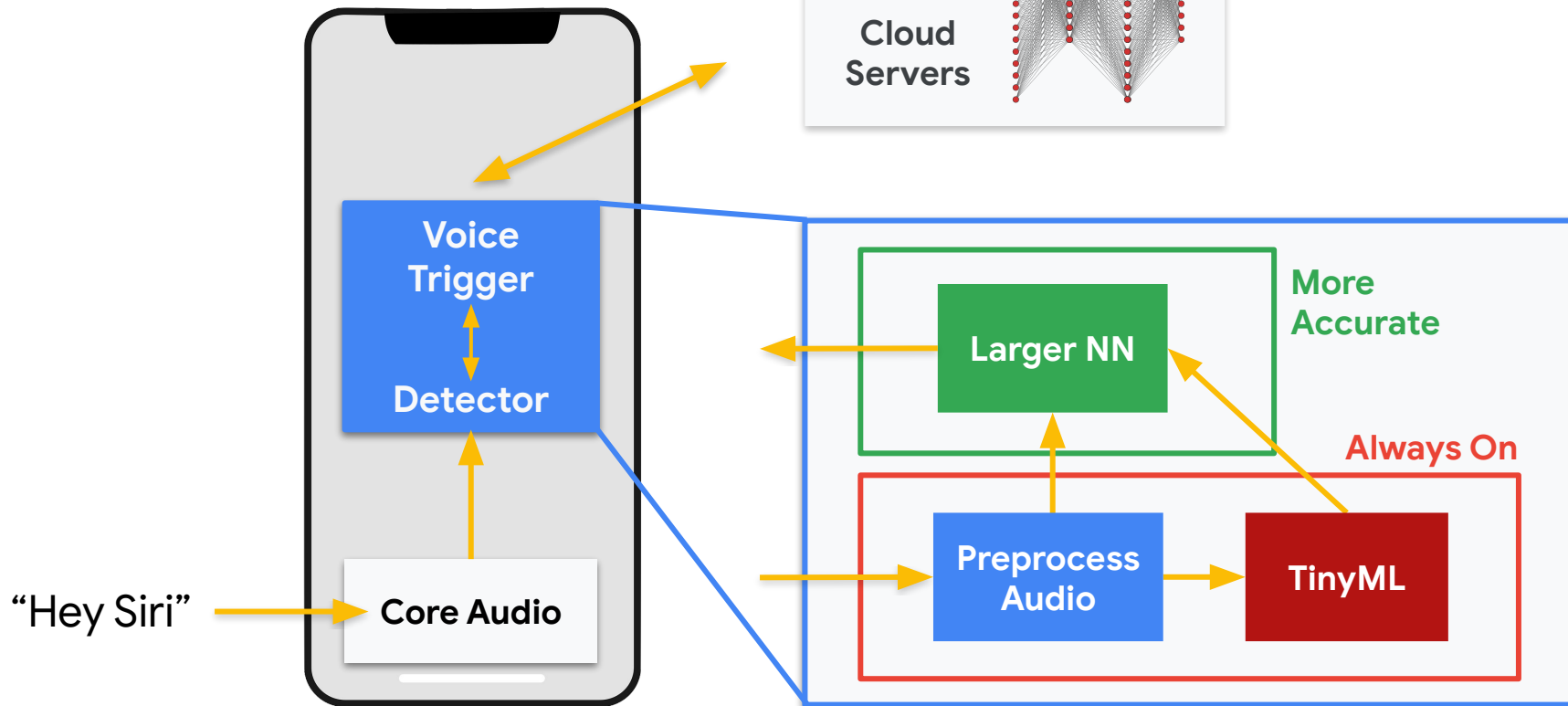
“Hey Siri”



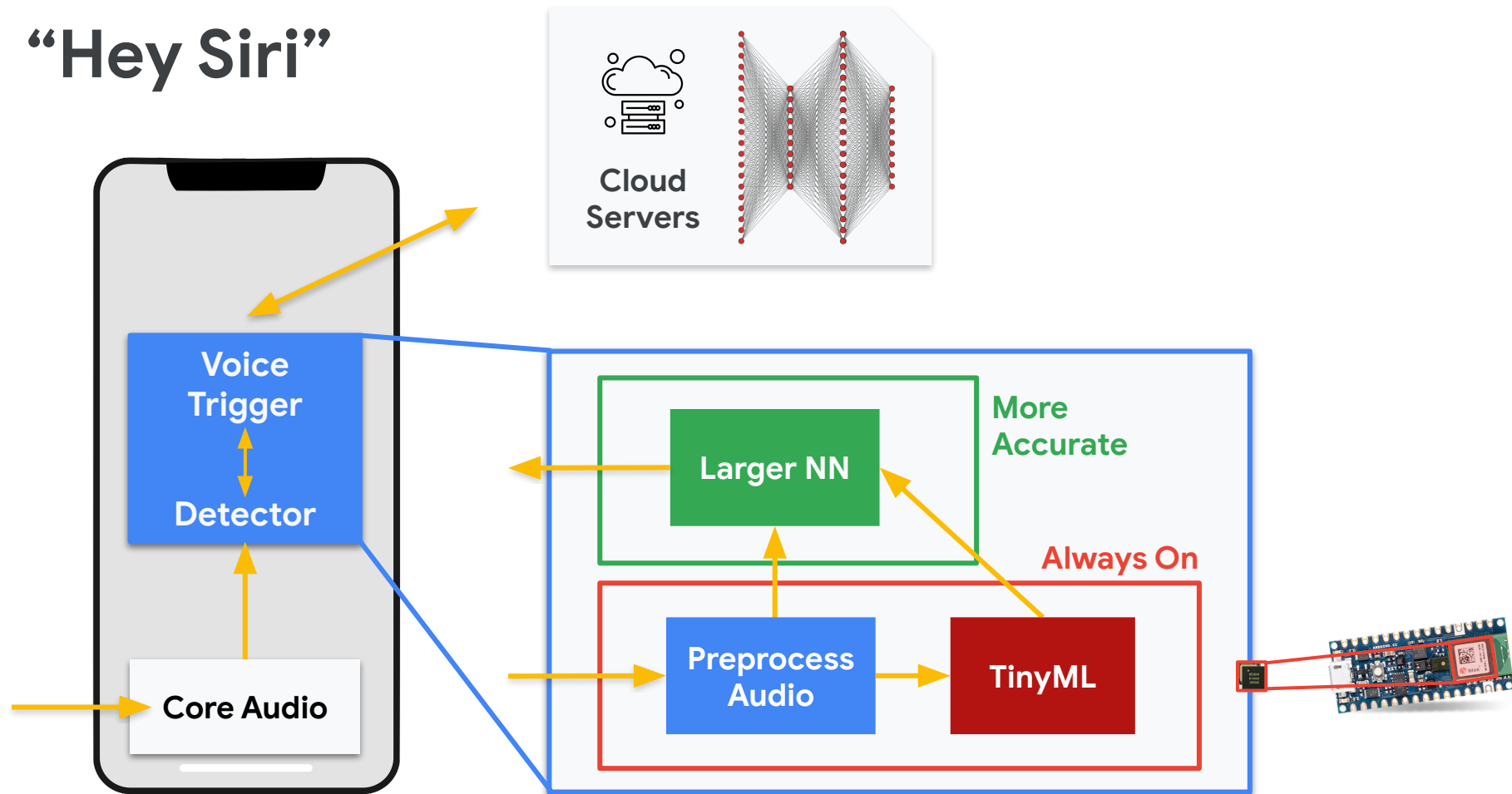
“Hey Siri”



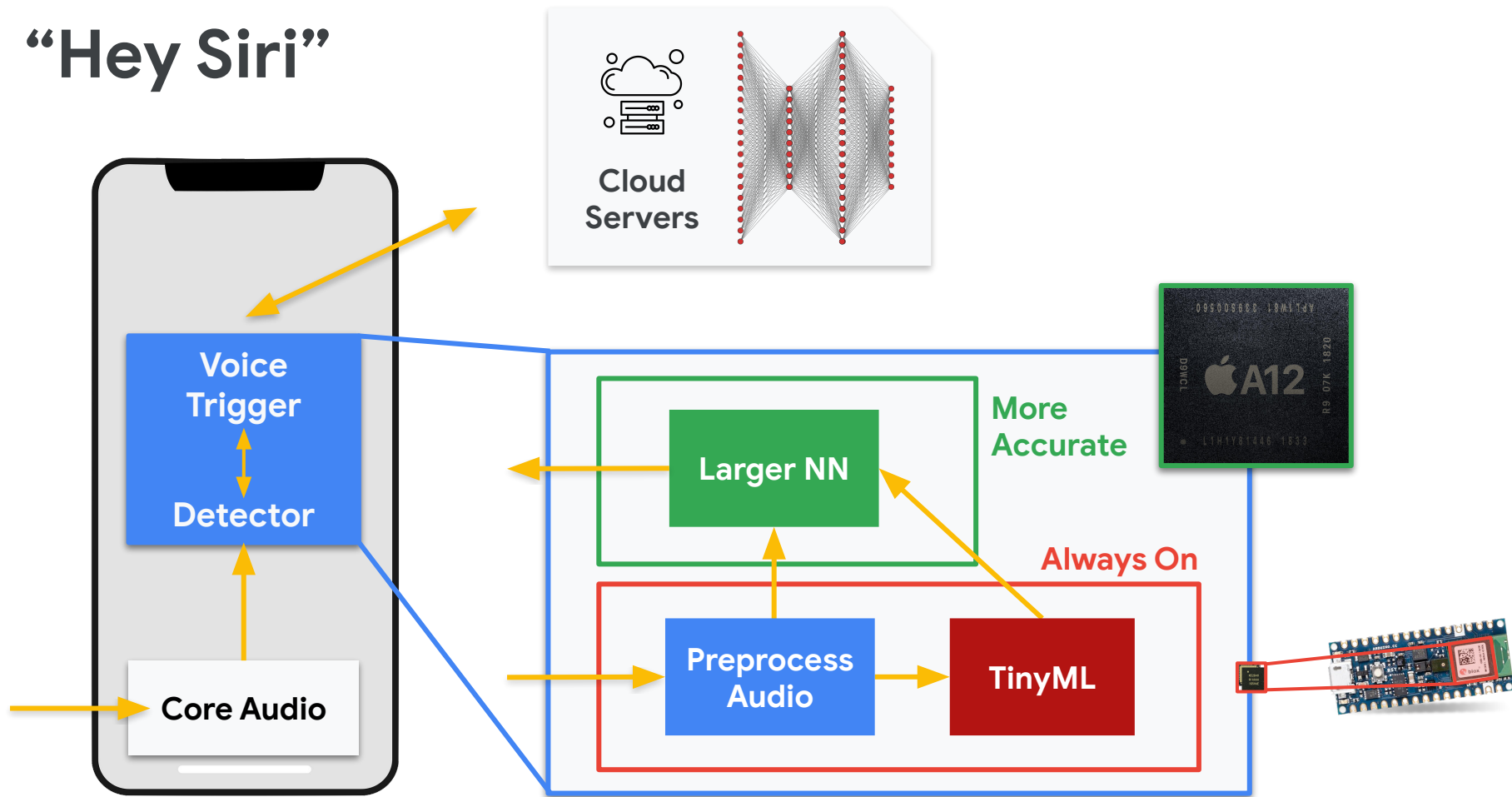
“Hey Siri”



“Hey Siri”

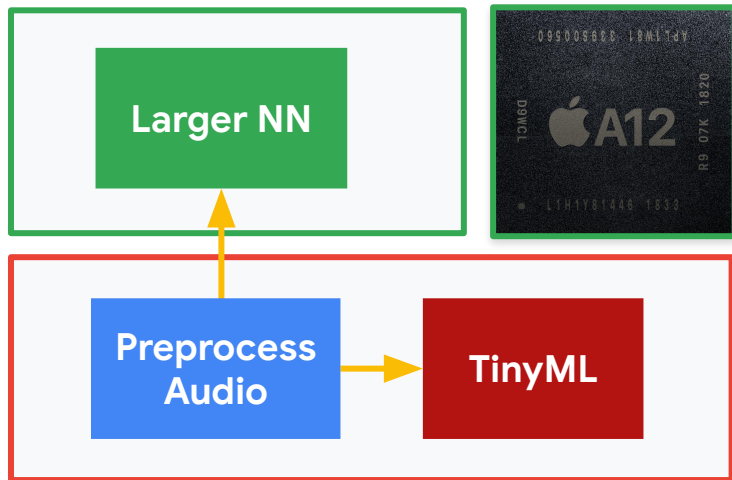


“Hey Siri”



“Hey Siri”

More Accurate



Always On

