

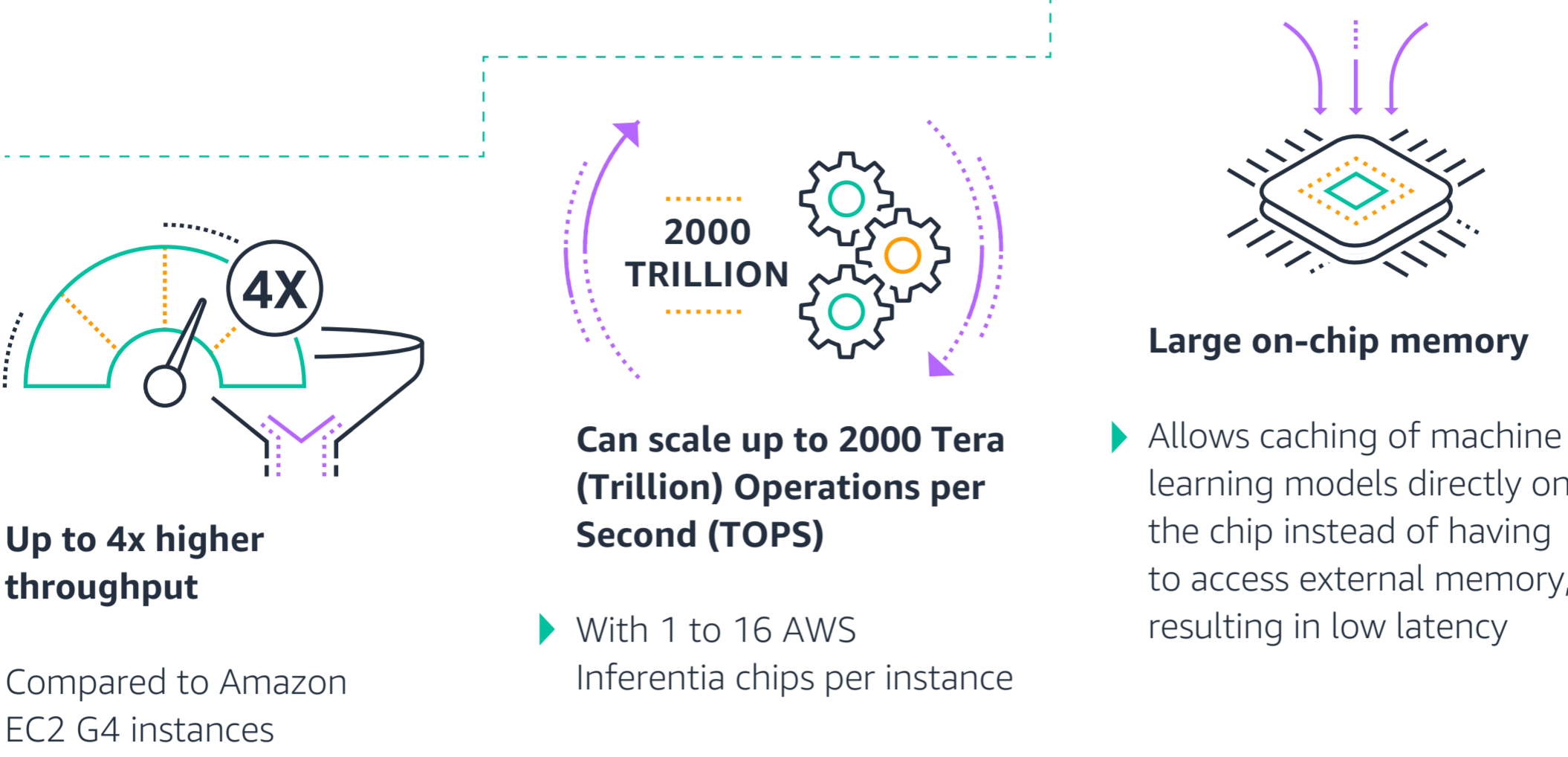
# Amazon EC2 Inf1 Instances: High Performance with the Lowest Cost Machine Learning Inference in the Cloud



With Amazon EC2 Inf1 instances powered by AWS Inferentia chips, you can optimize the deployment of your machine learning applications with high throughput, low latency, at the lowest cost per inference in the cloud.

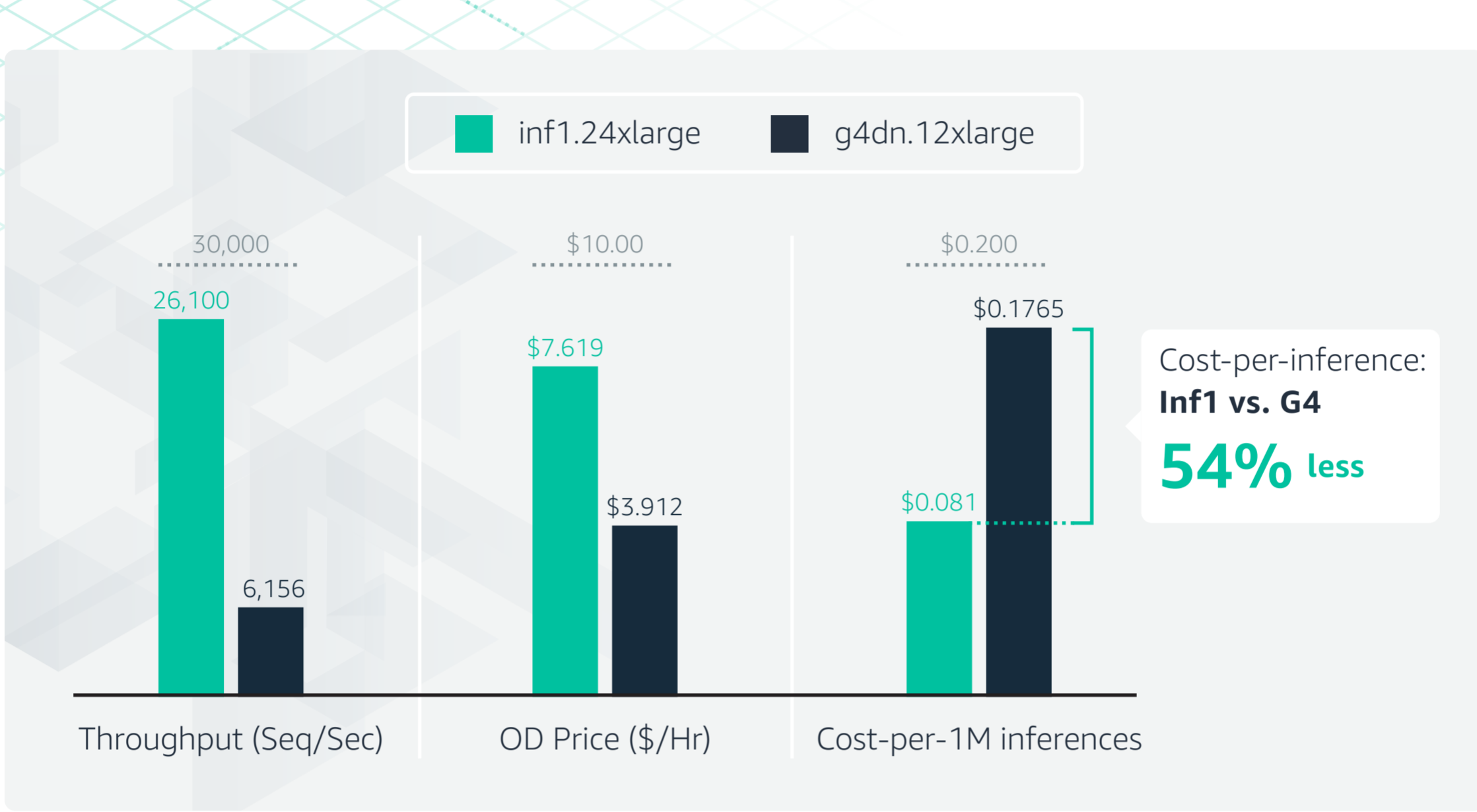
## Achieve optimized throughput and latency

High throughput and low latency mean you can achieve faster processing without compromise.



## Enable the lowest cost machine learning inference in the cloud

The cost of deploying a machine learning model can have a significant impact on budgets. Inf1 instances outperform other instances with the lowest cost per inference in the cloud.

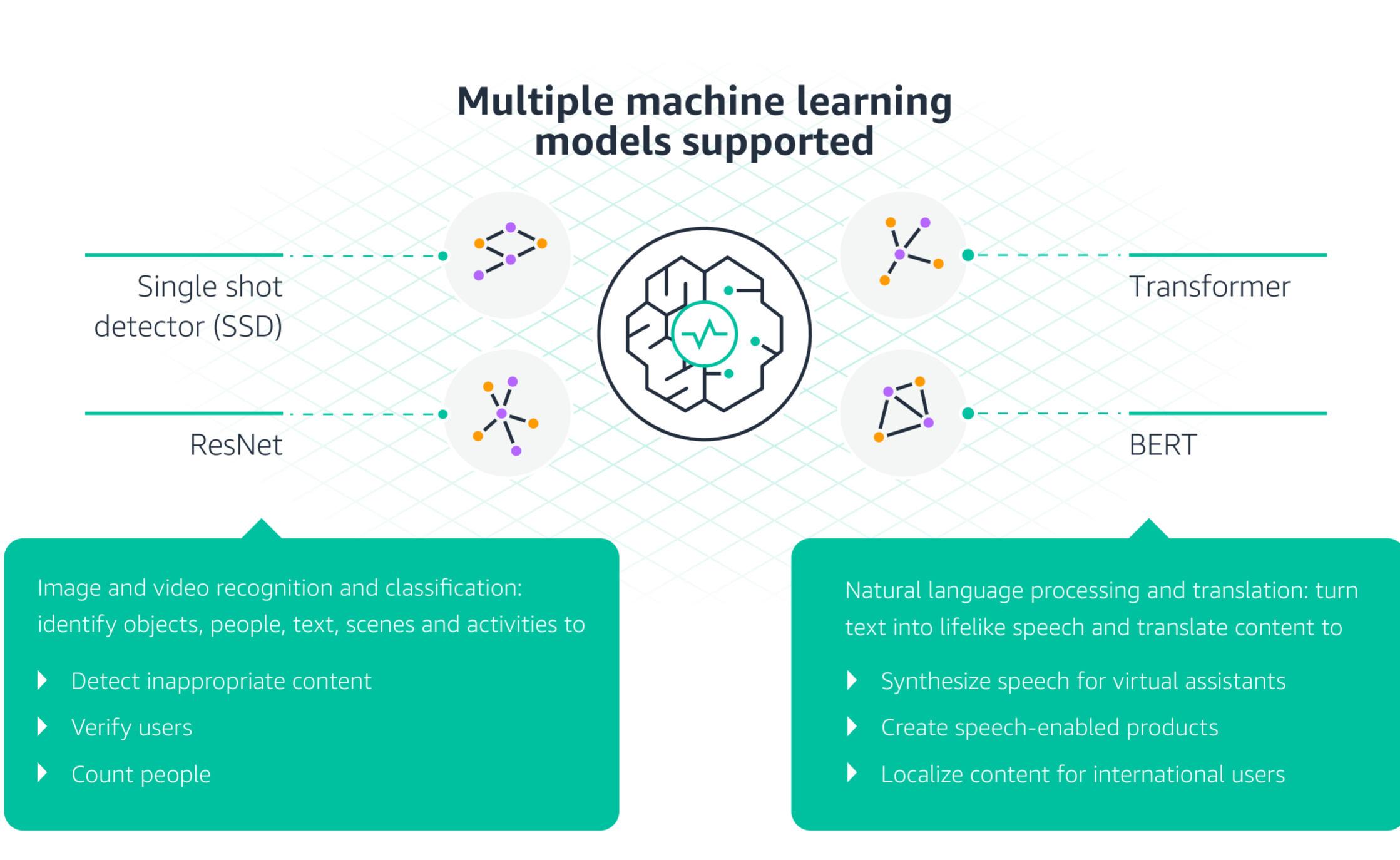


Amazon EC2 Inf1 instances deliver the lowest cost machine learning inference in the cloud



## Choose a flexible and easy-to-use solution

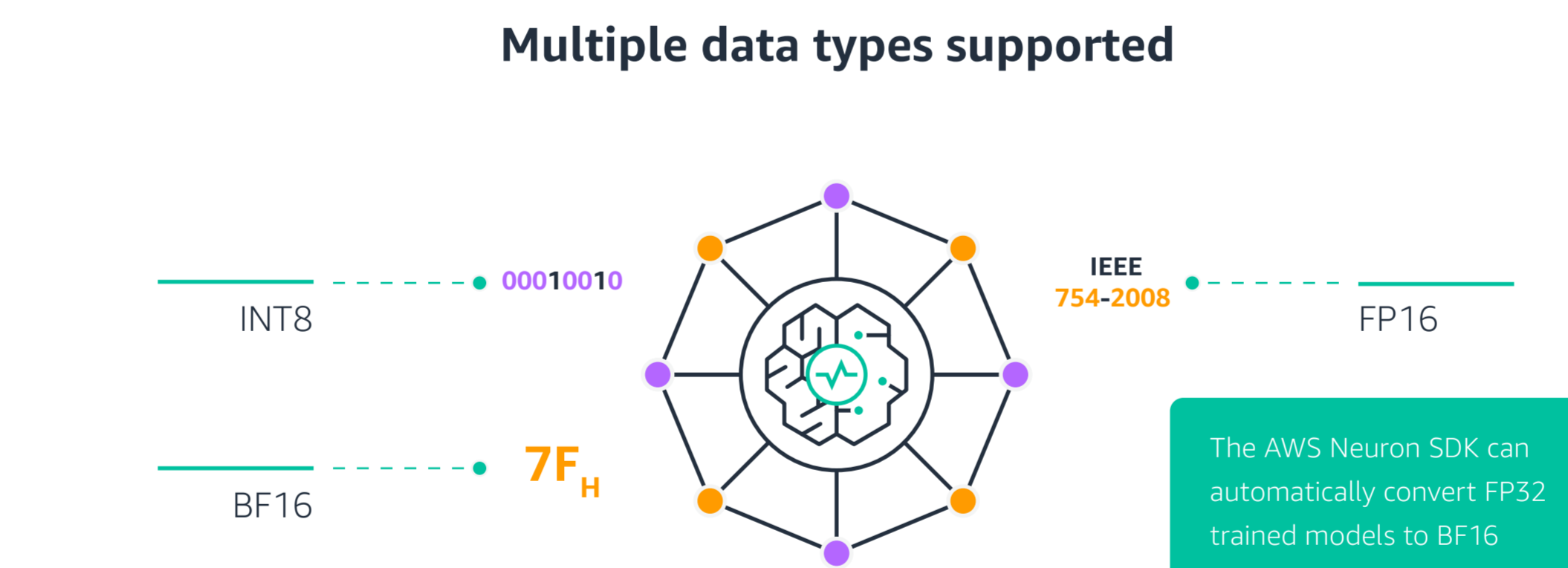
Inf1 instances support multiple machine learning models and data types, requiring few code changes to support models trained on the most popular frameworks.



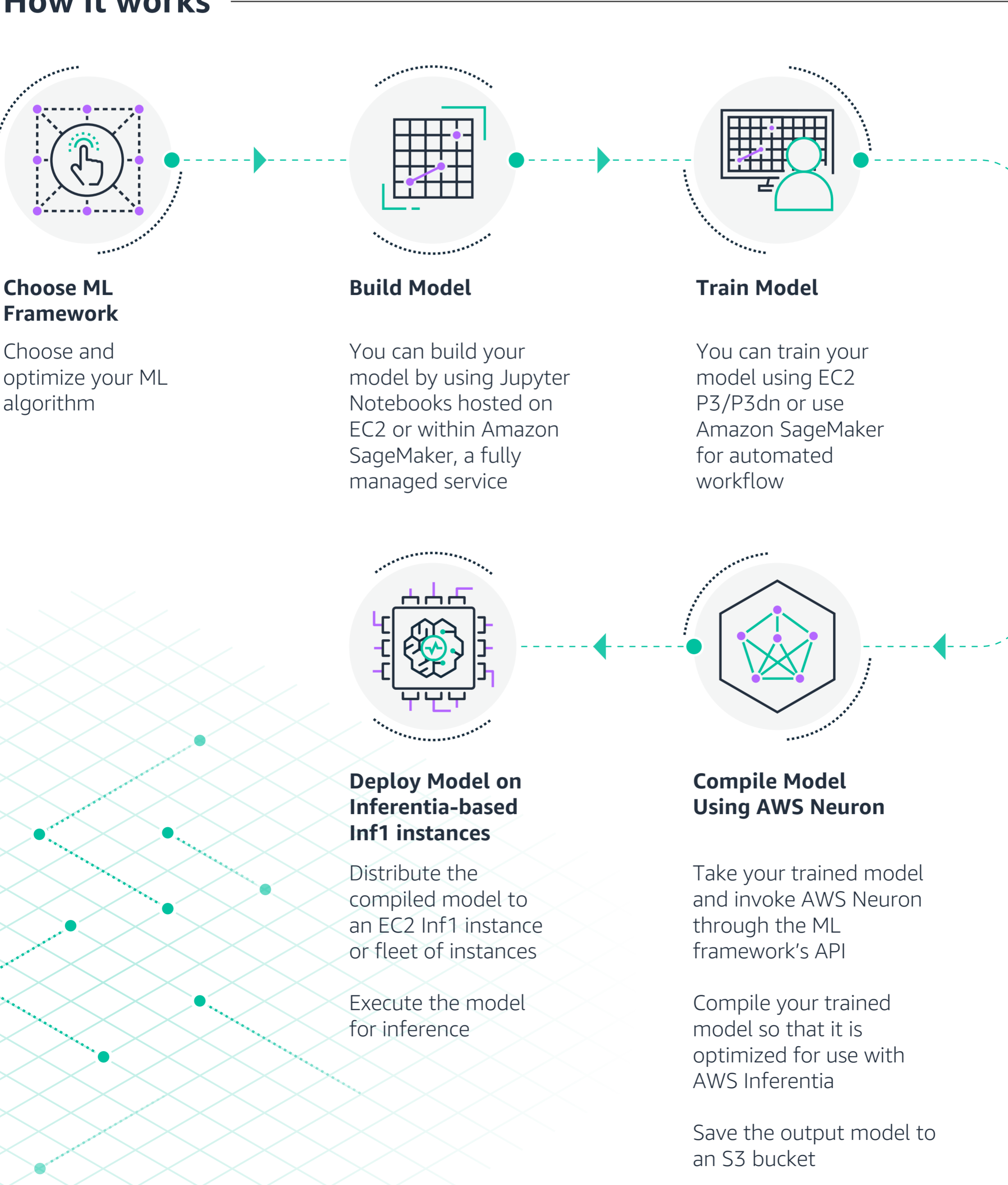
## Supports widely-used frameworks with few, if any, code changes



## Multiple data types supported



## How it works



With Amazon EC2 Inf1 instances, you can run a variety of large-scale ML inference applications at high throughput, low latency, at the lowest cost in the cloud.

Learn more at <https://aws.amazon.com/ec2/instance-types/inf1/>