

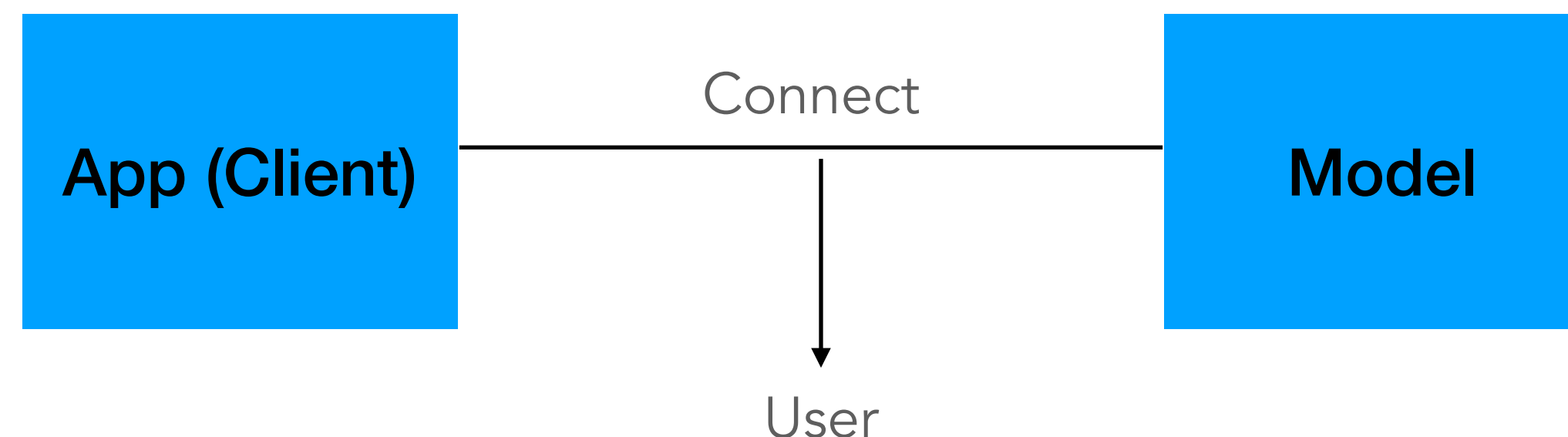
Mobile Nets

Enabling Computer Vision on the Web

Mobile Nets

Enabling Computer Vision on the Web

Let's assume we need to serve a (CV) model
in some sort of application!



Let's assume we do image classification in a web application.

➡ your app is going to be a Web App running on a browser (client).

Model on Server

- Inference is done on a server
- The client is thin i.e. only for I/O

Model on Client

- The server provides a model
- The client is thick, i.e., does inference as well

Mobile Nets

Bringing Computer Vision to the Web

Key Questions

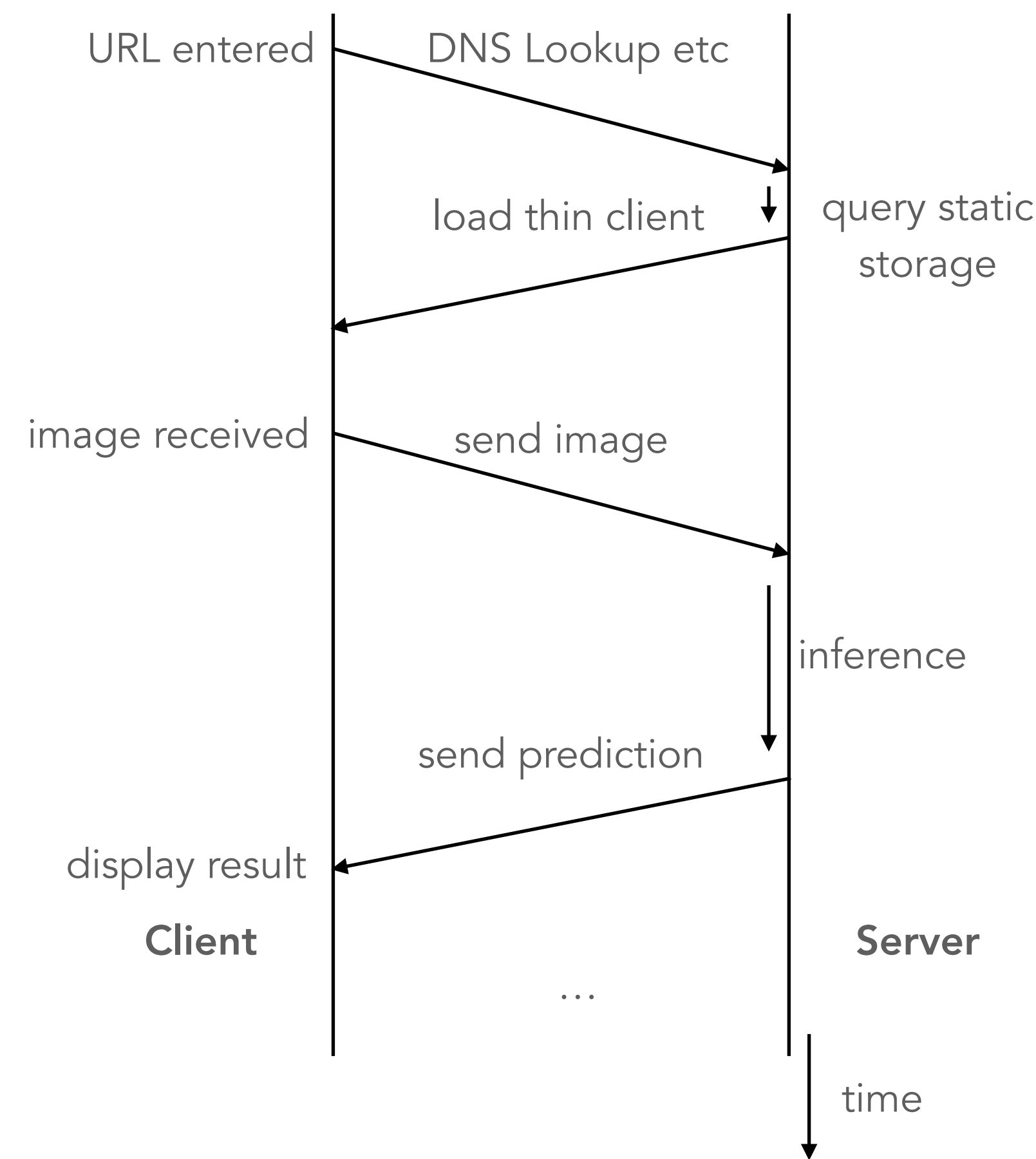
1. Why should I run a model on the browser?
2. How do Mobile Nets work and how are they helpful?
3. How **do** I run a model on the browser? (DEMO)

Why should I run a model on
the browser ?

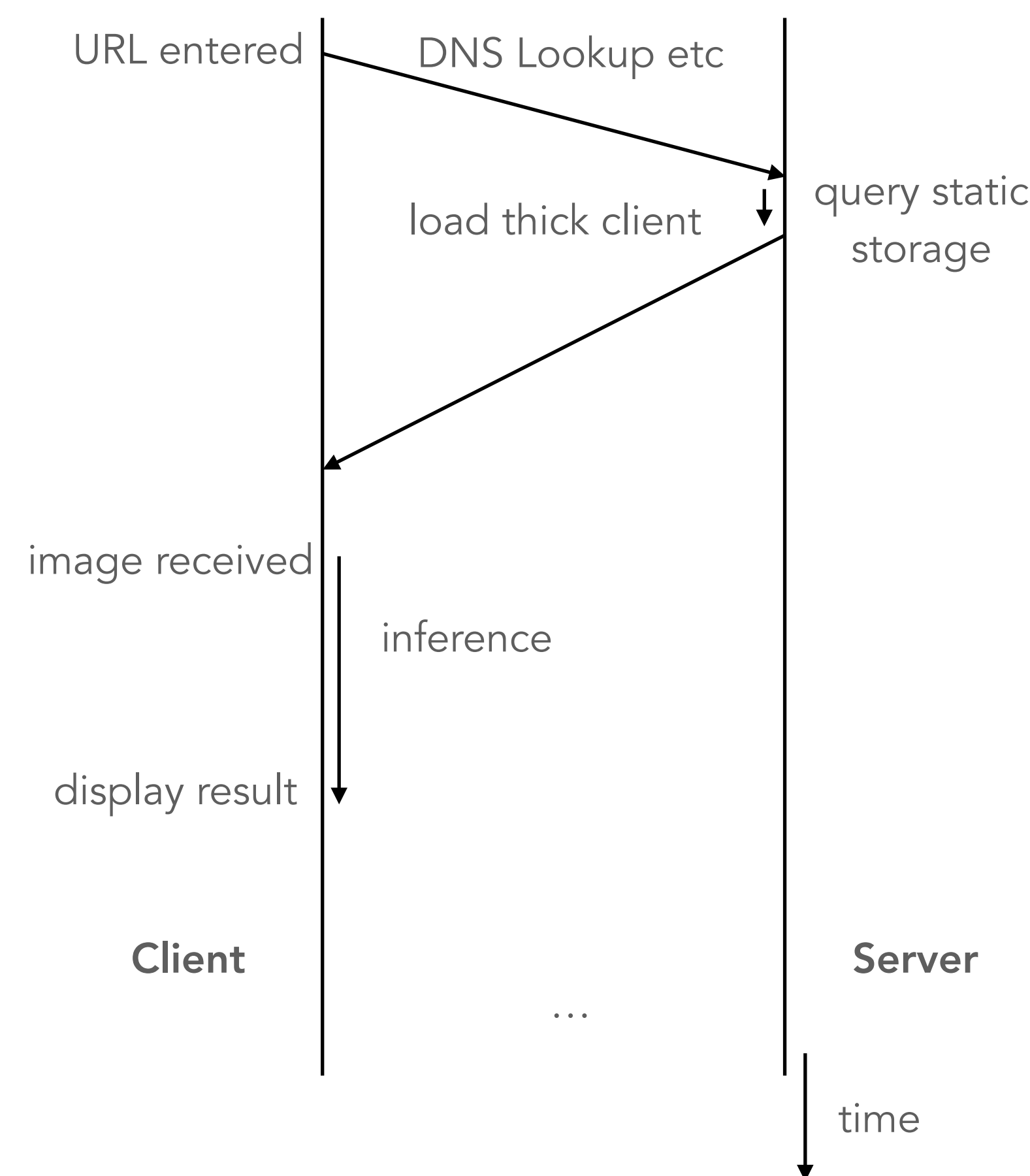
Why should I run a model on the browser ?

Pros/Cons for a browser-based application

Server-based Model



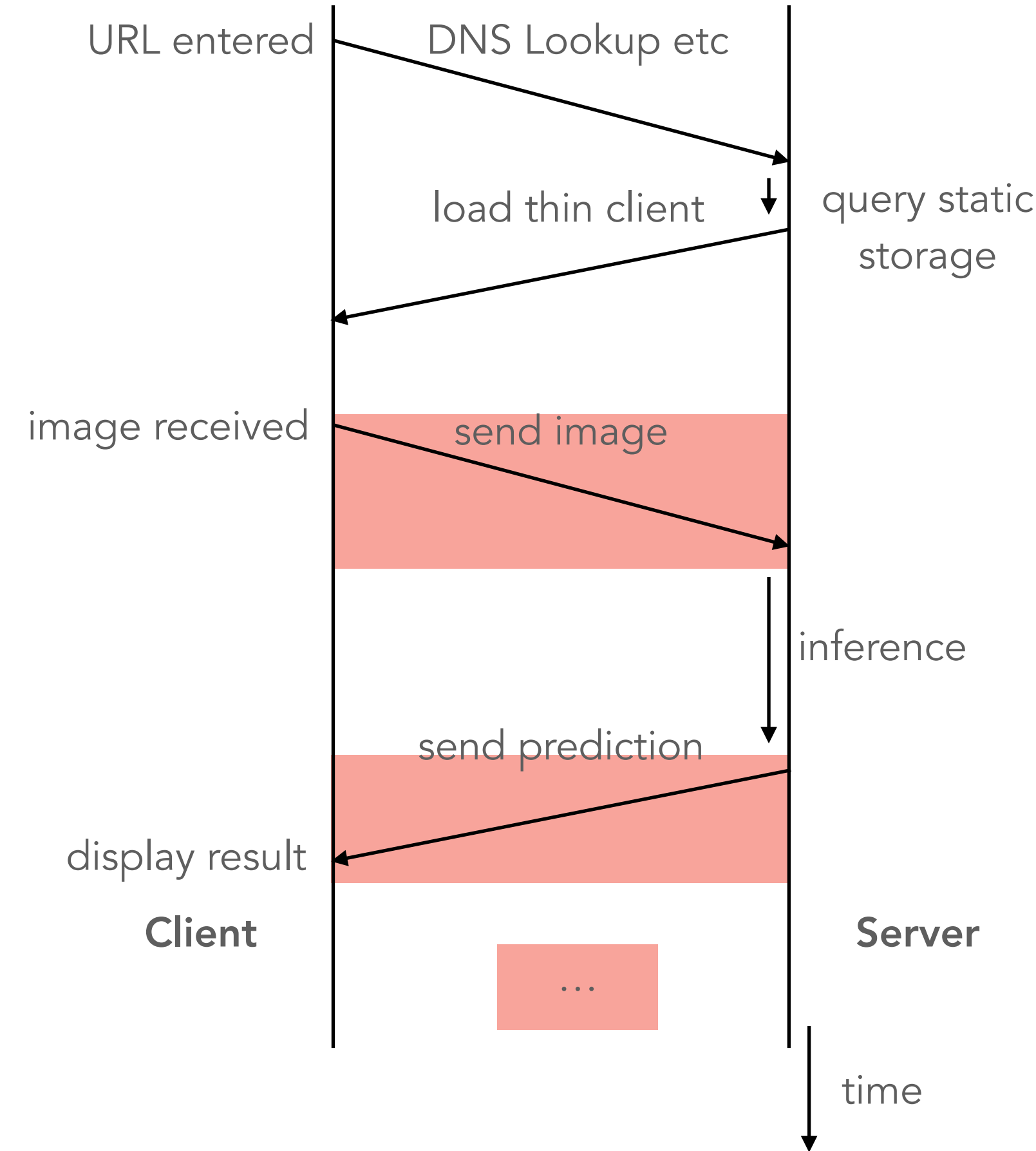
Client-based Model



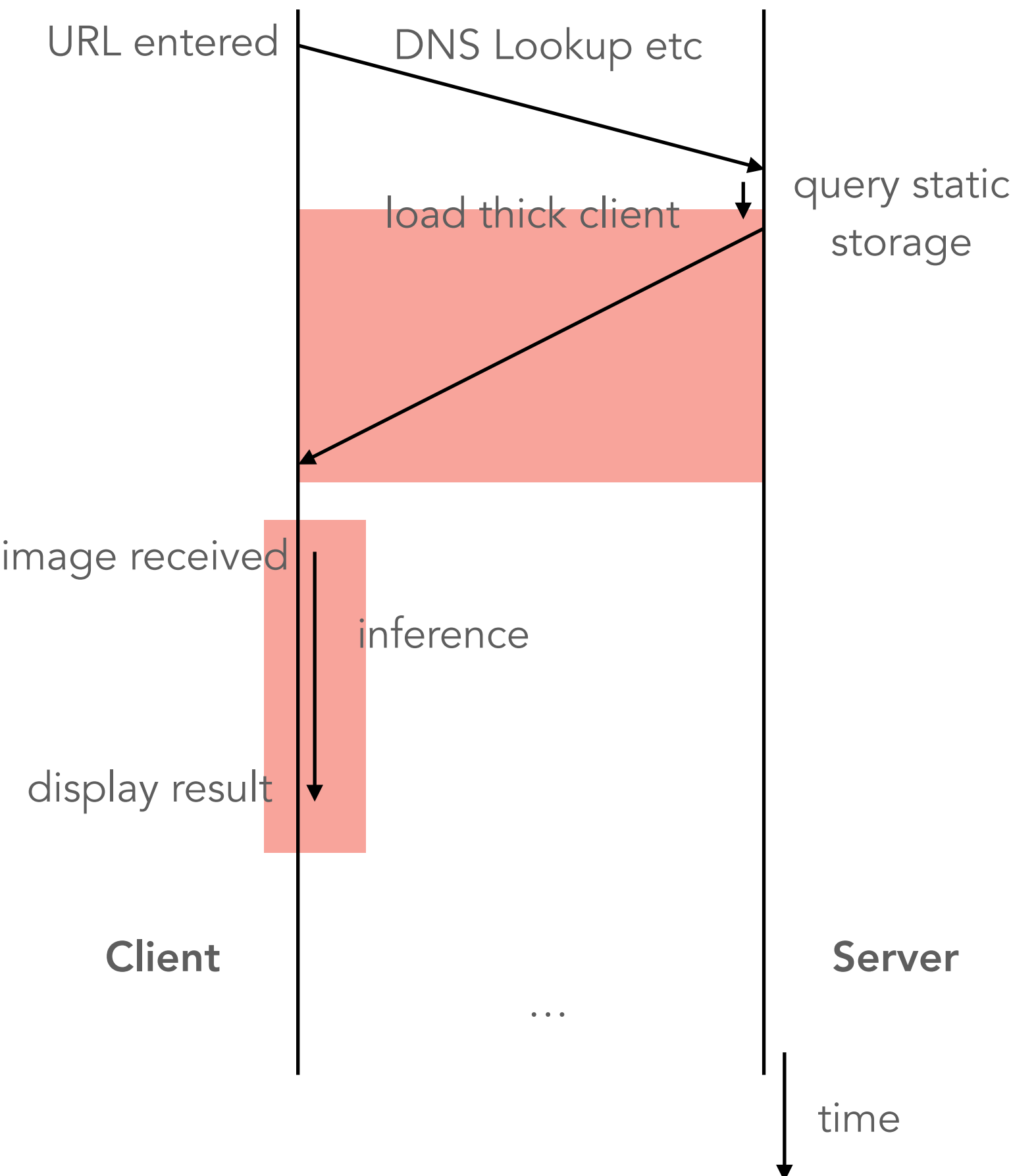
Why should I run a model on the browser ?

Pros/Cons for a browser-based application

Server-based Model



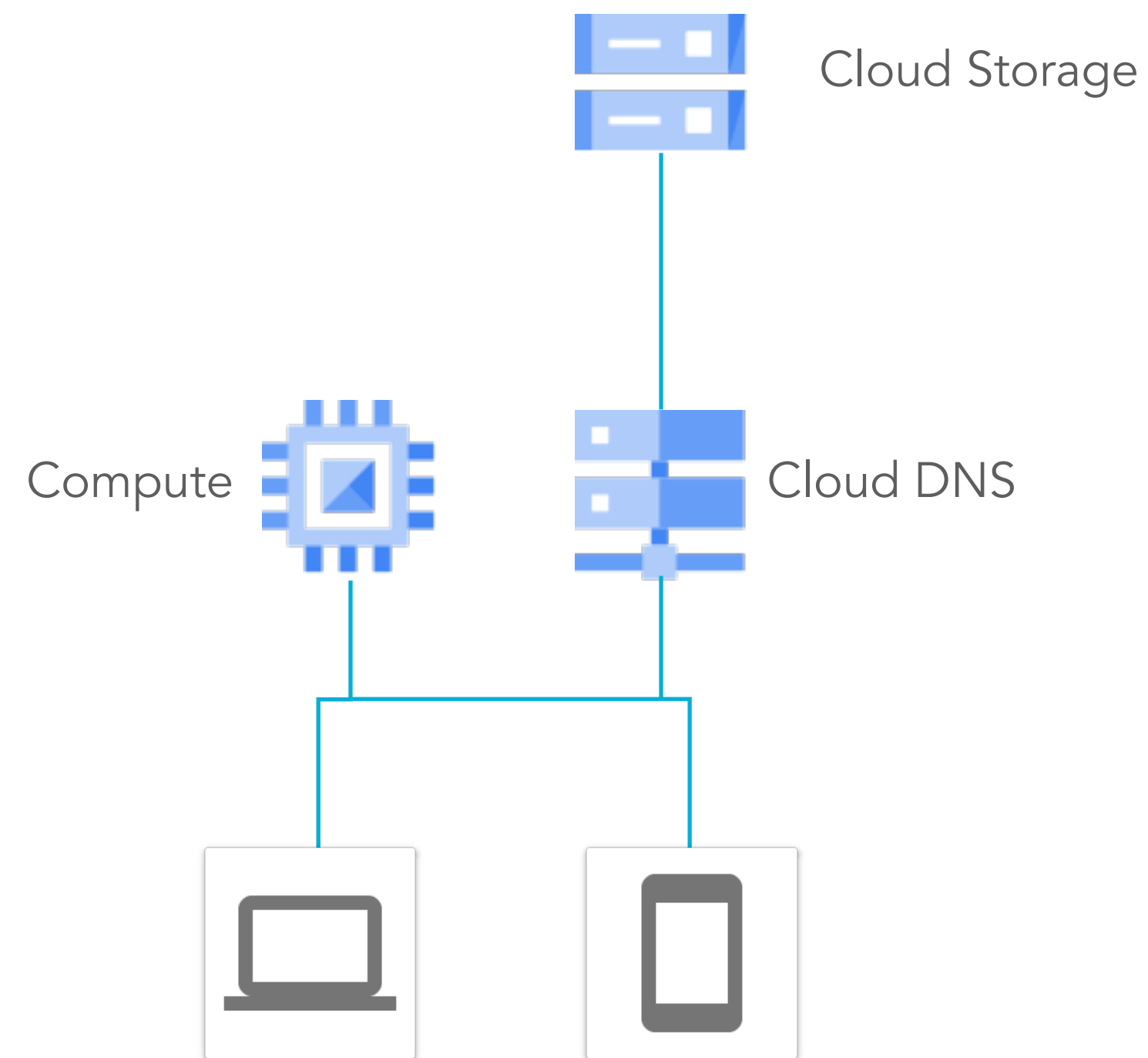
Client-based Model



Why should I run a model on the browser ?

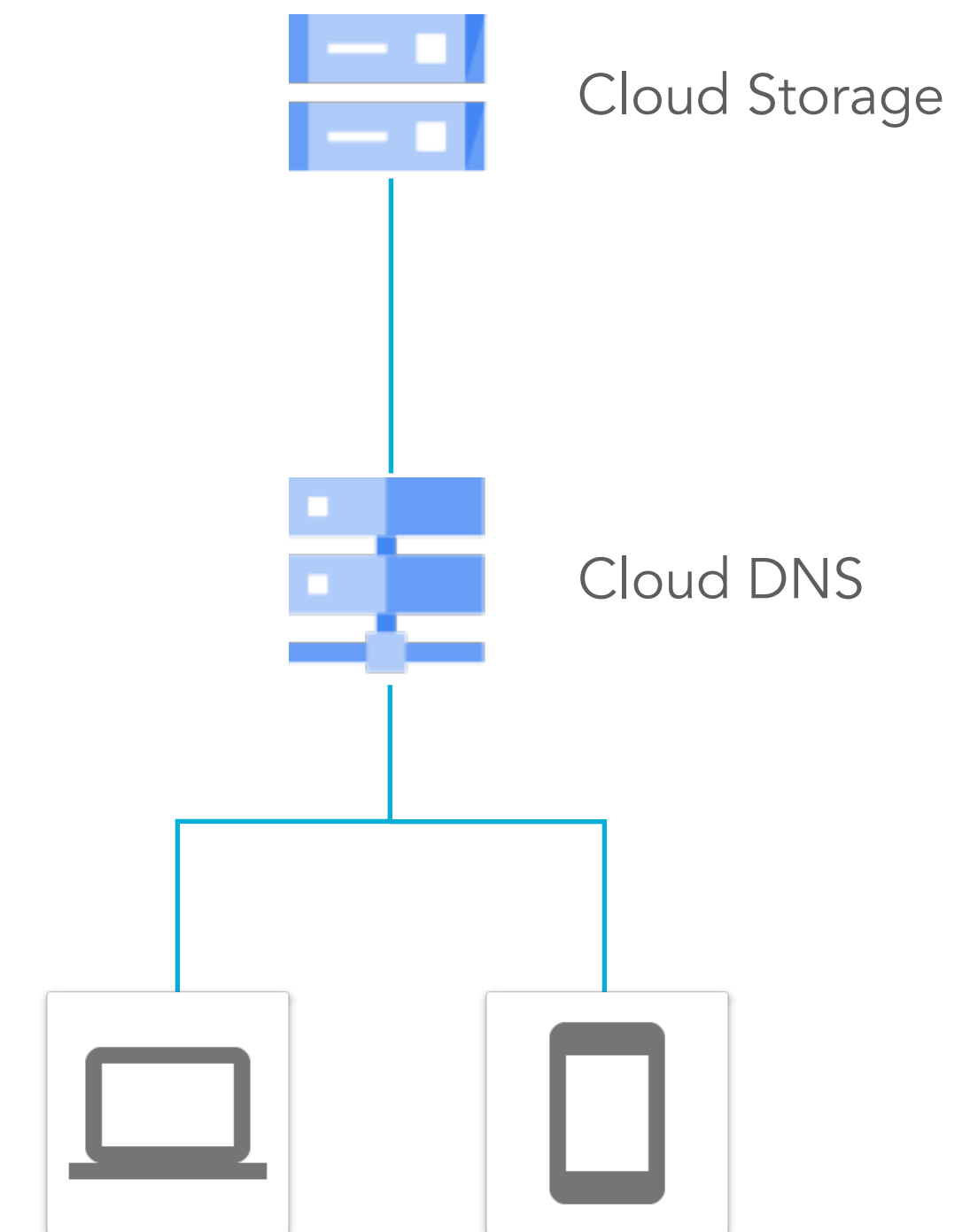
Pros/Cons for a browser-based application

Server-based Model



10 Users

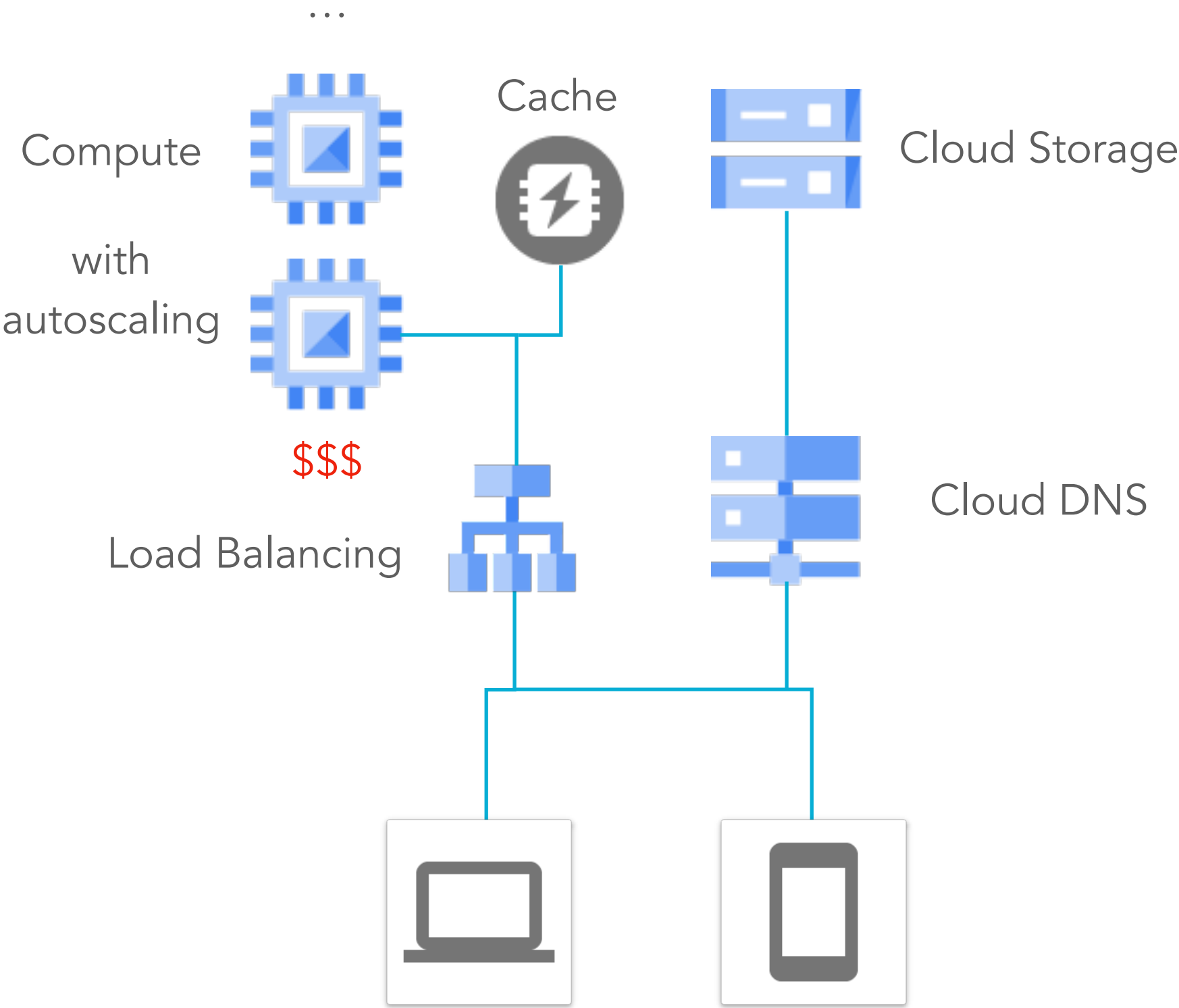
Client-based Model



Why should I run a model on the browser ?

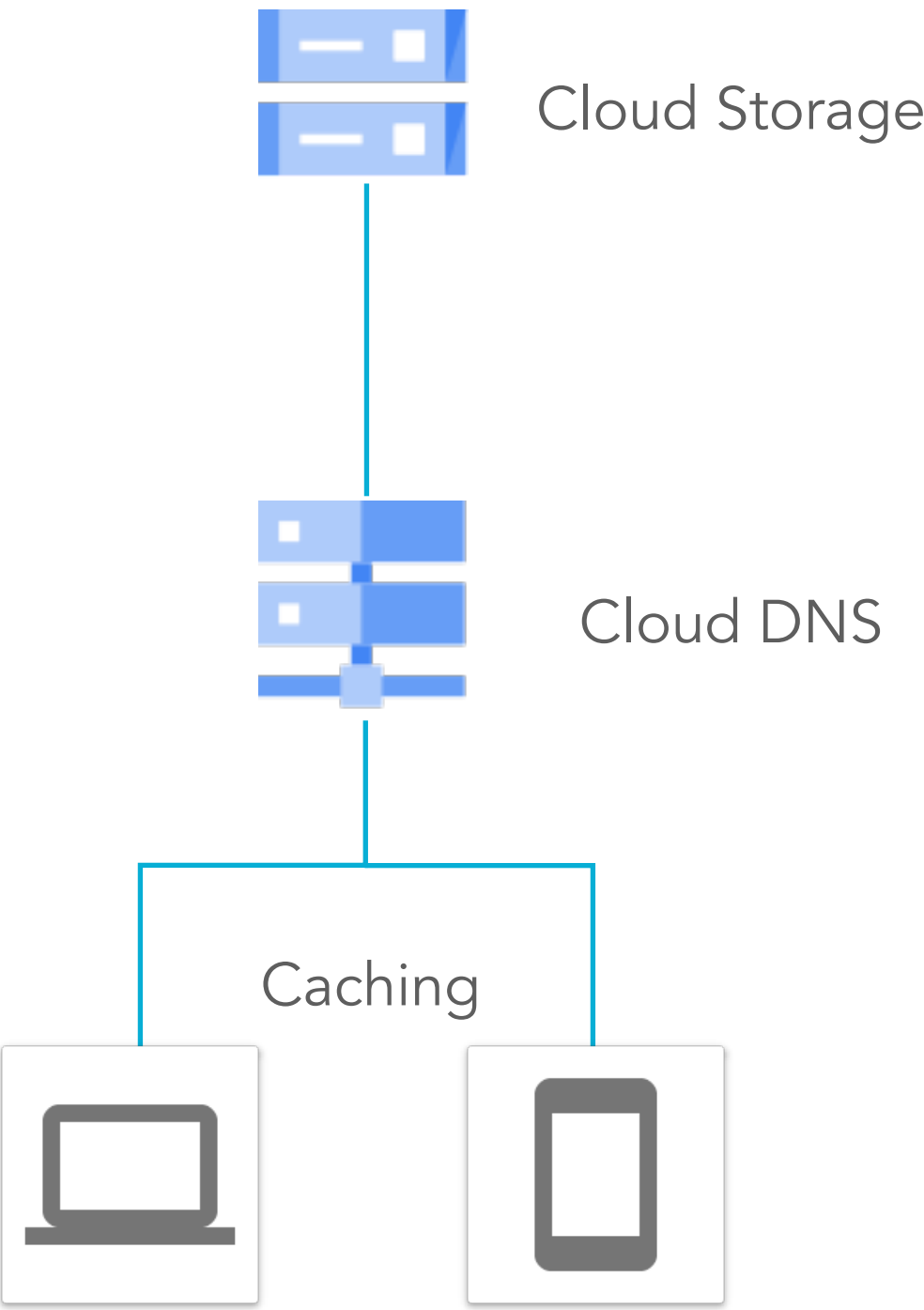
Pros/Cons for a browser-based application

Server-based Model



10000 Users

Client-based Model



Why should I run a model on the browser ?

Pros/Cons for a browser-based application

Server-based Model

+	-
unlimited memory and compute**	high sustained network traffic
	privacy
	needs reliable connection
	cost scales with no. predictions
	architecture complexity scales

Client-based Model

+	-
low latency	slow initial load
privacy	memory limitations
cost hardly scales	compute limitations
simple architecture	

Why should I run a model on the browser ?

Pros/Cons for a browser-based application

Server-based Model

+	-
unlimited memory and compute**	high sustained network traffic
	privacy
	needs reliable connection
	cost scales with no. predictions
	architecture complexity scales

Client-based Model

+	-
low latency	slow initial load
privacy	memory limitations
cost hardly scales	compute limitations
simple architecture	

Why should I run a model on
the browser ?

**If size and complexity of the model allow it,
why NOT run it on the browser?**

How do Mobile Nets work and how are they helpful?

How do mobile nets work and how are they helpful?

What problem do Mobile Nets address?

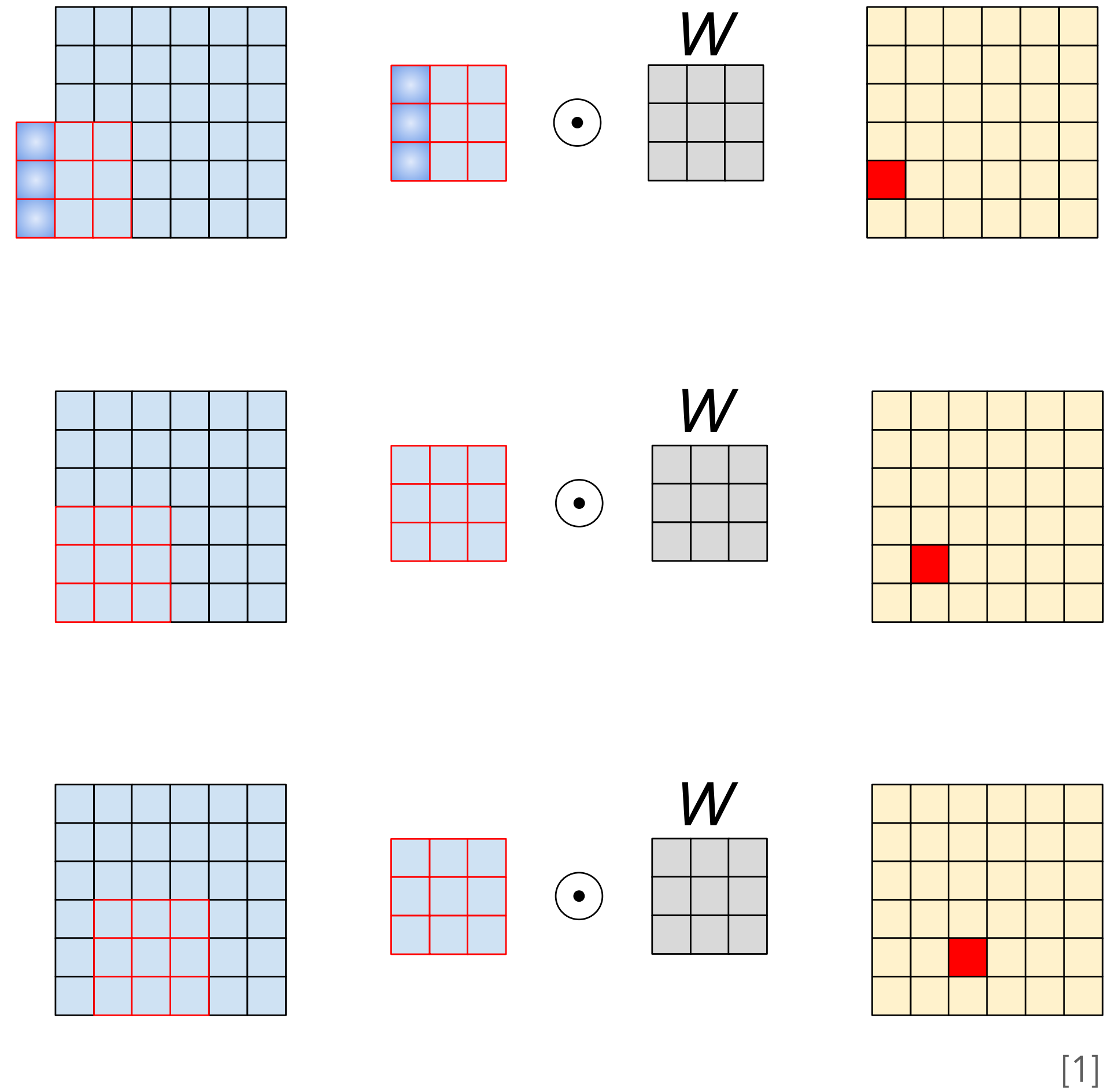
- ▶ CV models are large and computationally complex
- ▶ This makes them hard to run on the client side.
- ▶ In “MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications” Howard et al. aim to develop SOA Classifiers which can run on the client
- ▶ They observed: Models make heavy use of convolutional operations

Mobile Nets use depth-wise separable convolutions to reduce size & complexity

Mobile Nets introduce model shrinking parameters to fine-tune size & complexity

How do mobile nets work and how are they helpful?

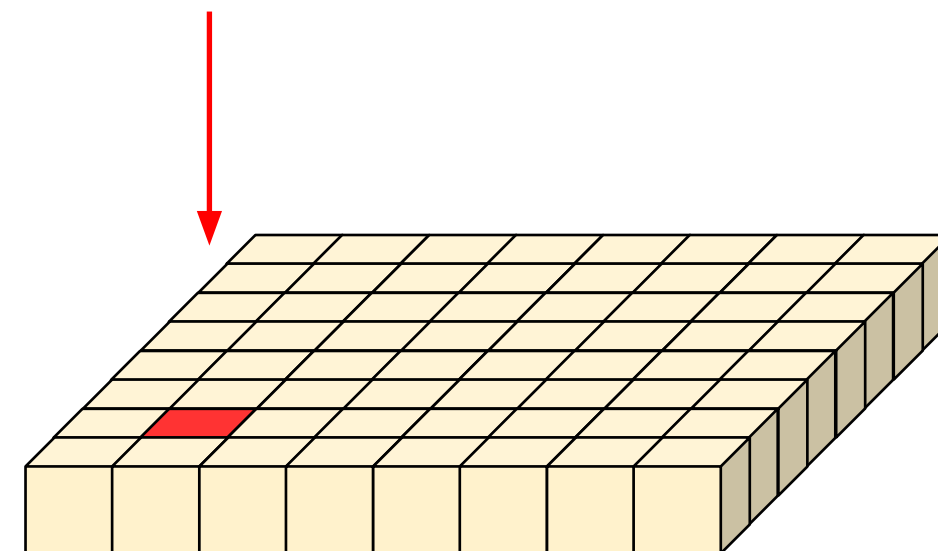
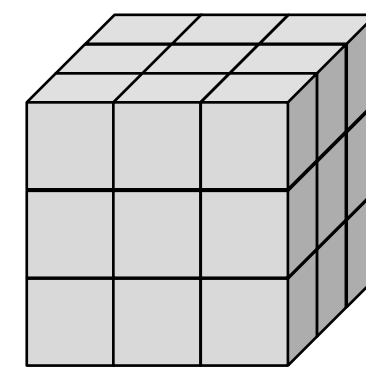
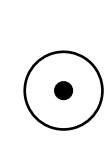
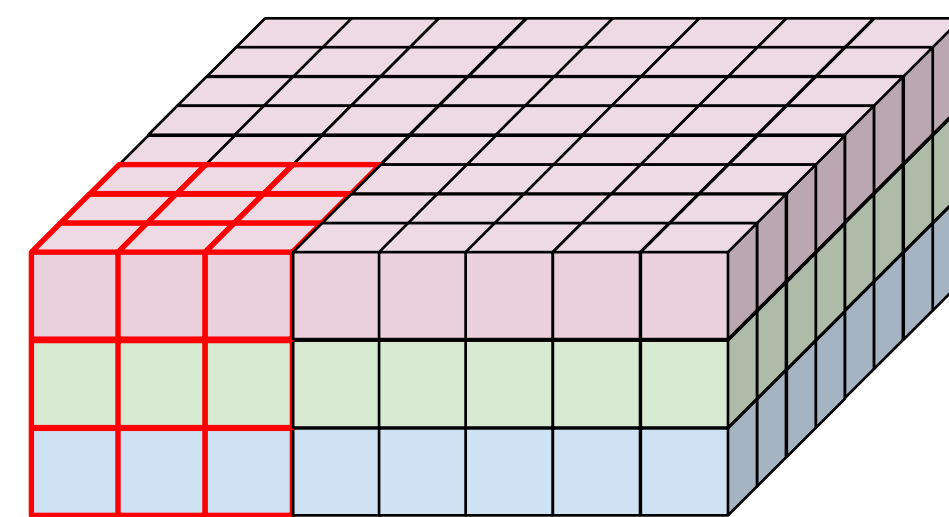
Depth-Wise Separable Convolutions



2D convolution

How do mobile nets work and how are they helpful?

Depth-Wise Separable Convolutions



[1]

3D convolution

Lets assume:

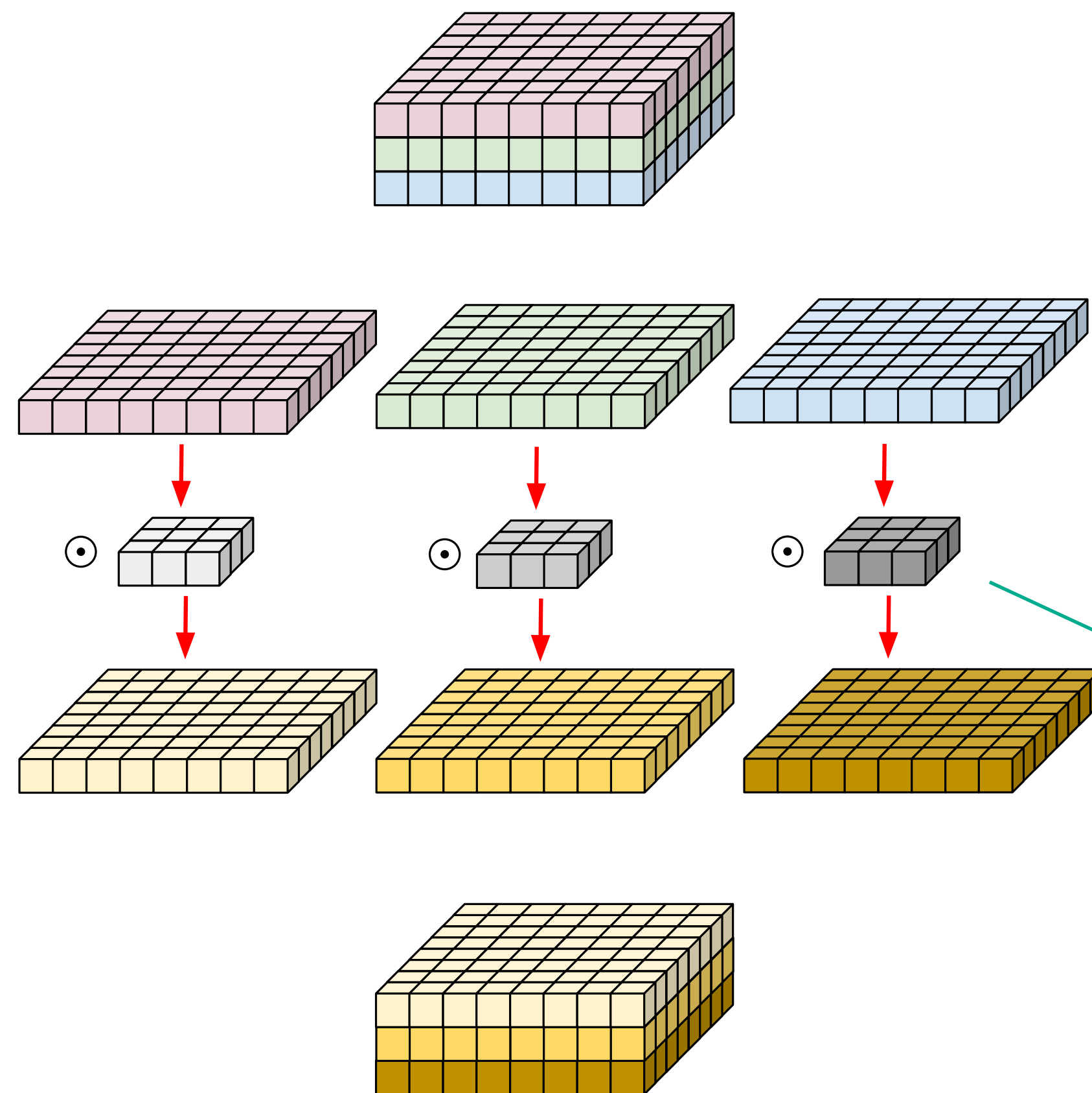
- Image size of 8x8
- padding of 1
- stride of 1
- 3 input channels
- kernel size 3x3
- N output channels

Parameters:

$$3 \times 3 \times 3 \times N = 27 \times N$$

How do mobile nets work and how are they helpful?

Depth-Wise Separable Convolutions



depth-wise convolution

Lets assume:

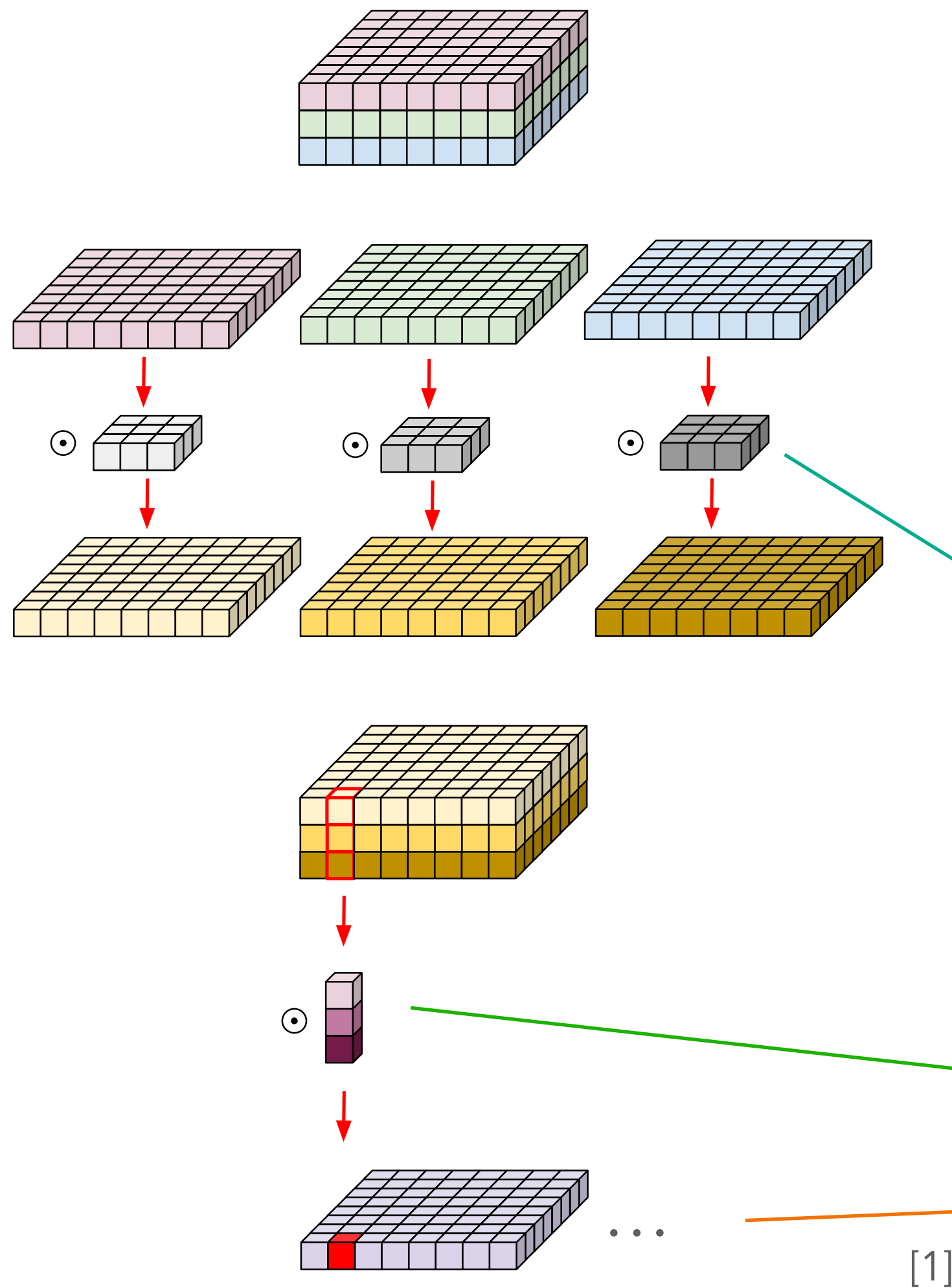
- Image size of 8x8
- padding of 1
- stride of 1
- 3 input channels
- kernel size 3x3
- N output channels

Parameters:

$$3 \times 3 \times 3 = 27$$

How do mobile nets work and how are they helpful?

Depth-Wise Separable Convolutions



Lets assume:

- Image size of 8x8
- padding of 1
- stride of 1
- 3 input channels
- kernel size 3x3
- N output channels

Parameters:

27

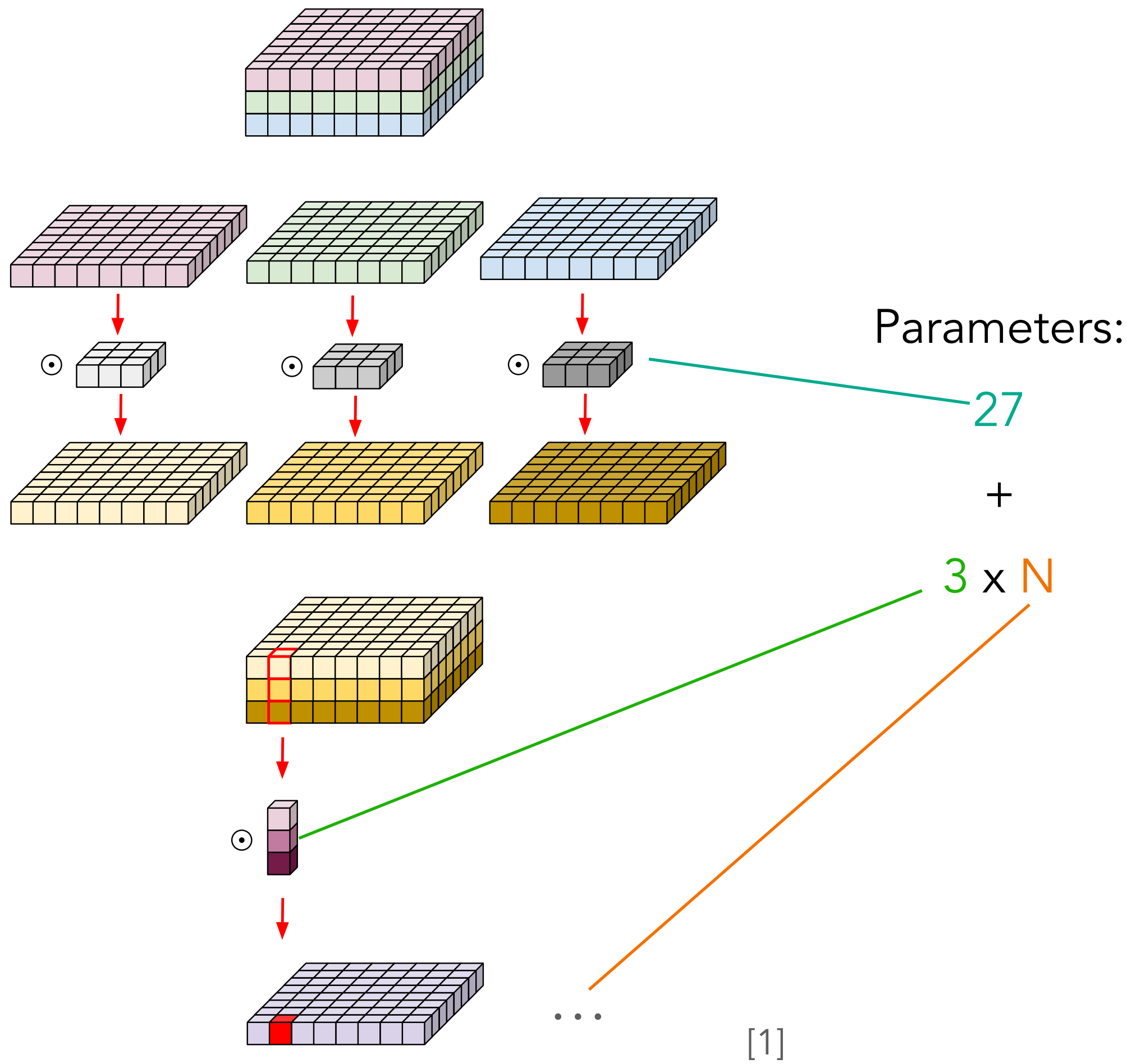
+

3 x N

depth-wise separable convolution

How do mobile nets work and how are they helpful?

Depth-Wise Separable Convolutions

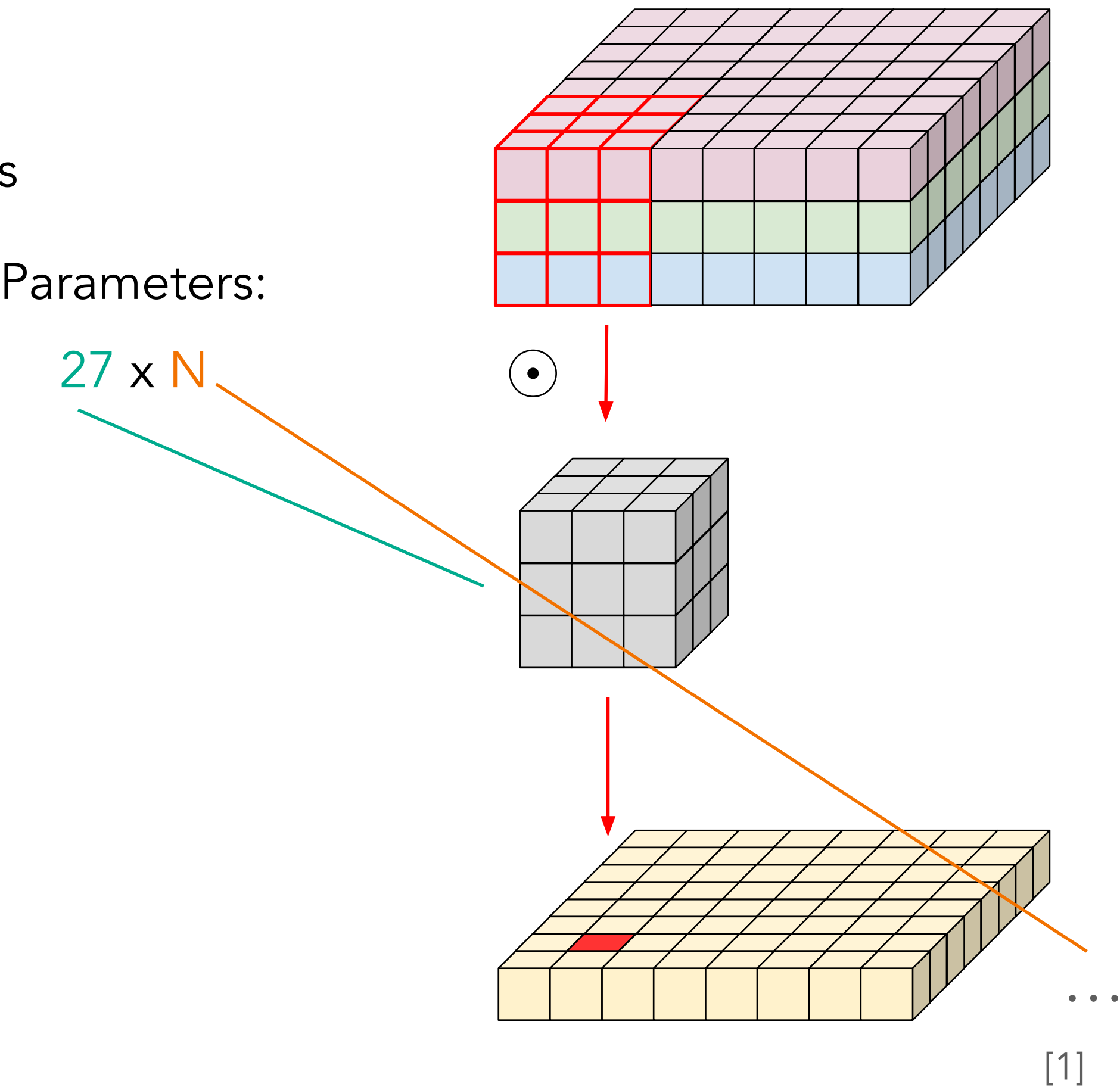


depth-wise separable convolution

- 3 input channels
- kernel size 3x3
- N output channels

Advantage

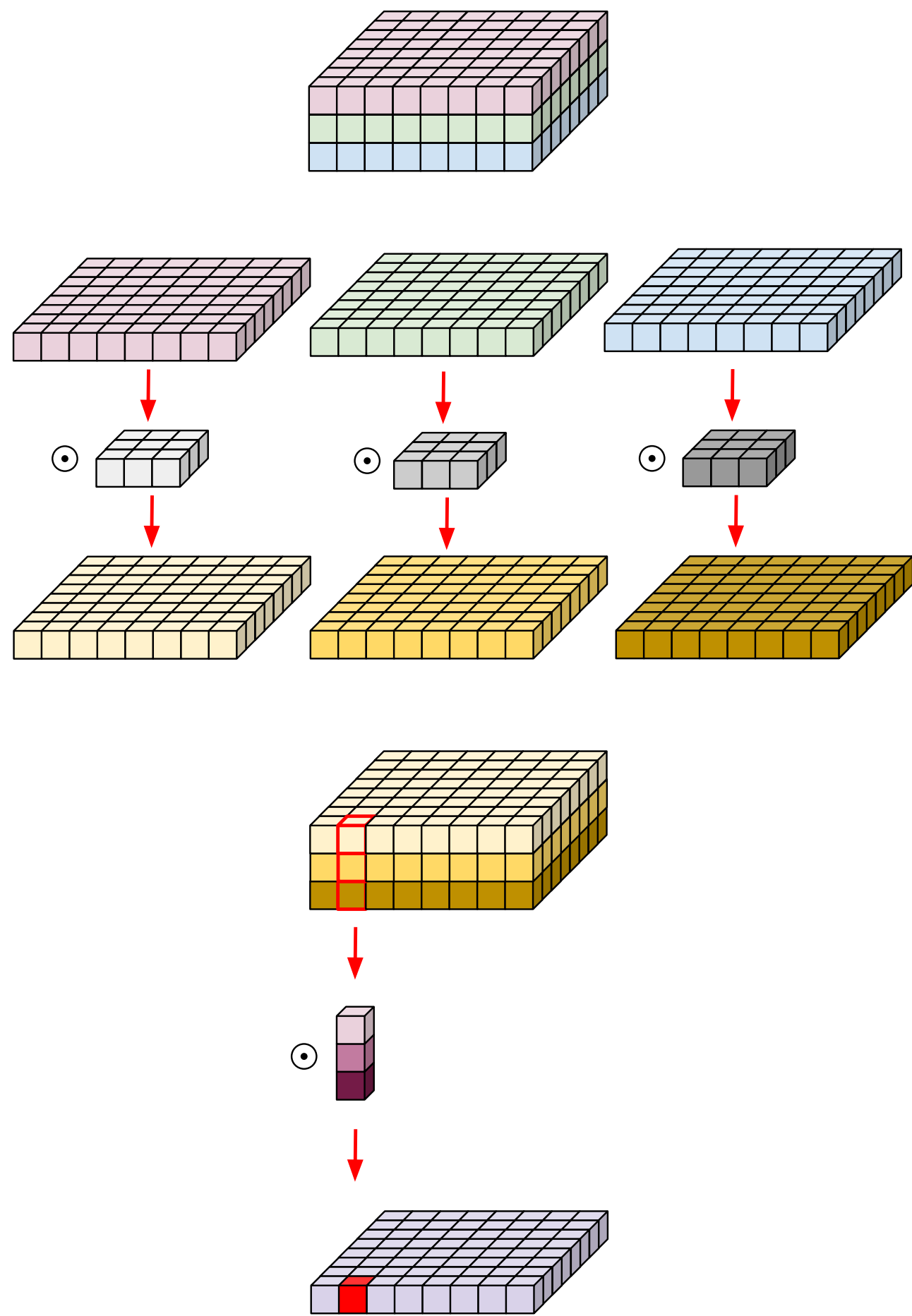
$$\frac{1}{N} + \frac{1}{9}$$



3D convolution

How do mobile nets work and how are they helpful?

Depth-Wise Separable Convolutions



depth-wise separable convolution

[1]

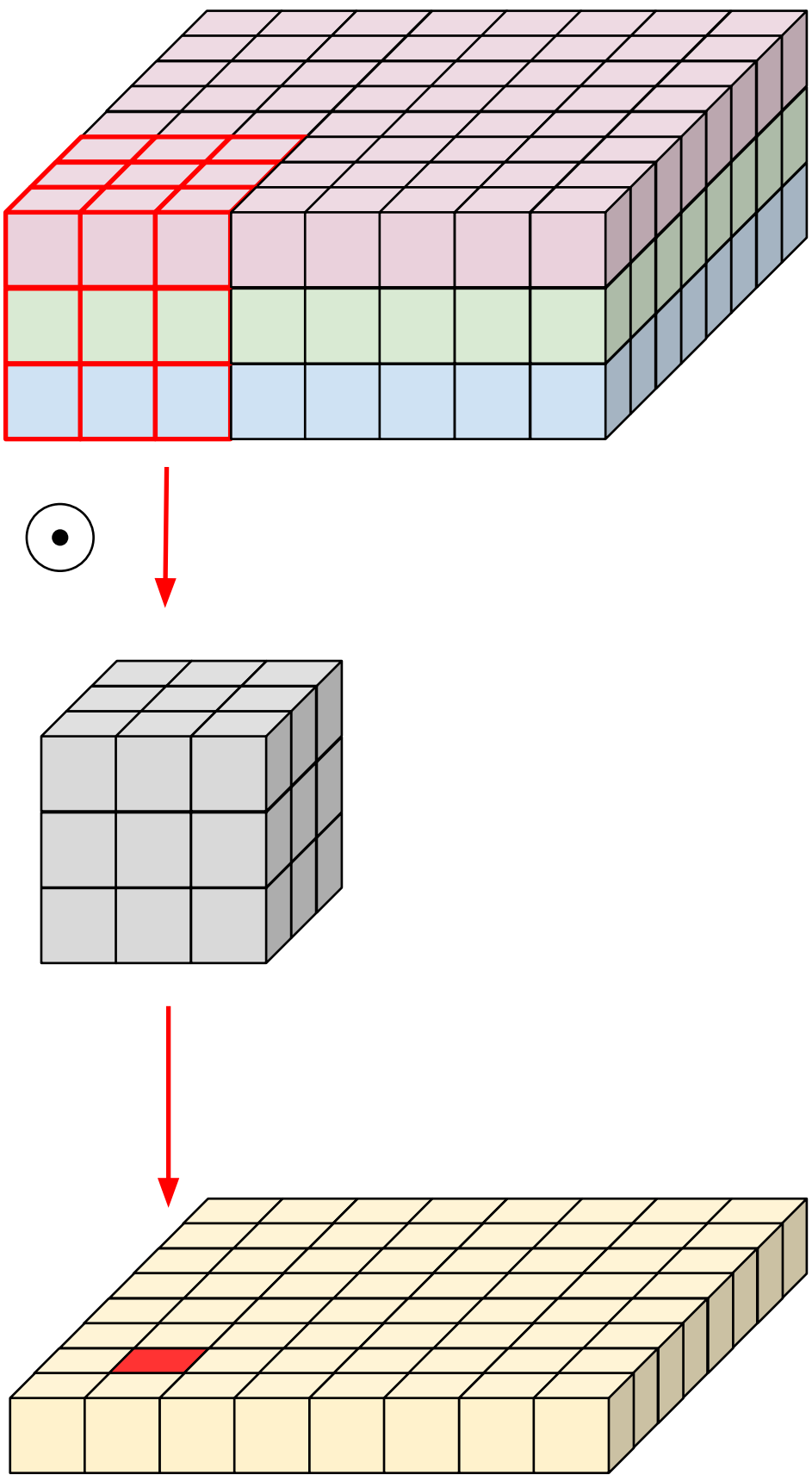
Limitation:
Our filters cannot mix input channels!

It seems to have little effect in practice:

Table 4. Depthwise Separable vs Full Convolution MobileNet

Model	ImageNet	Million	Million
	Accuracy	Mult-Adds	Parameters
Conv MobileNet	71.7%	4866	29.3
MobileNet	70.6%	569	4.2

[2]



3D convolution

[1]

How do mobile nets work and how are they helpful?

Model Scaling Parameters

- ▶ In addition to depth-wise separable convolutions, two model scaling parameters are introduced.
- ▶ α , set in $(0,1]$ scales the number of input and output channels in each layer.
- ▶ Reduces the number of parameters and reduces computational cost by roughly α^2
- ▶ ρ set in $(0,1]$, scales the resolution of the images and subsequently the resolution of each internal representation (between the layers).
- ▶ Reduces computational cost by ρ^2

Layer/Modification	Million Mult-Adds	Million Parameters
Convolution	462	2.36
Depthwise Separable Conv	52.3	0.27
$\alpha = 0.75$	29.6	0.15
$\rho = 0.714$	15.1	0.15

[2]

How do Mobile Nets work and how are they helpful?

Mobile nets are competitive, efficient CNNs, which can easily be scaled down in terms of size and complexity.

How do I run a model on the browser?

DEMO

Any Questions?



[mx-e/cv-seminar-deployment-showcase](https://github.com/mx-e/cv-seminar-deployment-showcase)

Sources

- [1] A really great blog post about depth-separable convolutions:
<https://eli.thegreenplace.net/2018/depthwise-separable-convolutions-for-machine-learning/>
- [2] The Mobile Nets paper:
<https://arxiv.org/pdf/1704.04861.pdf>

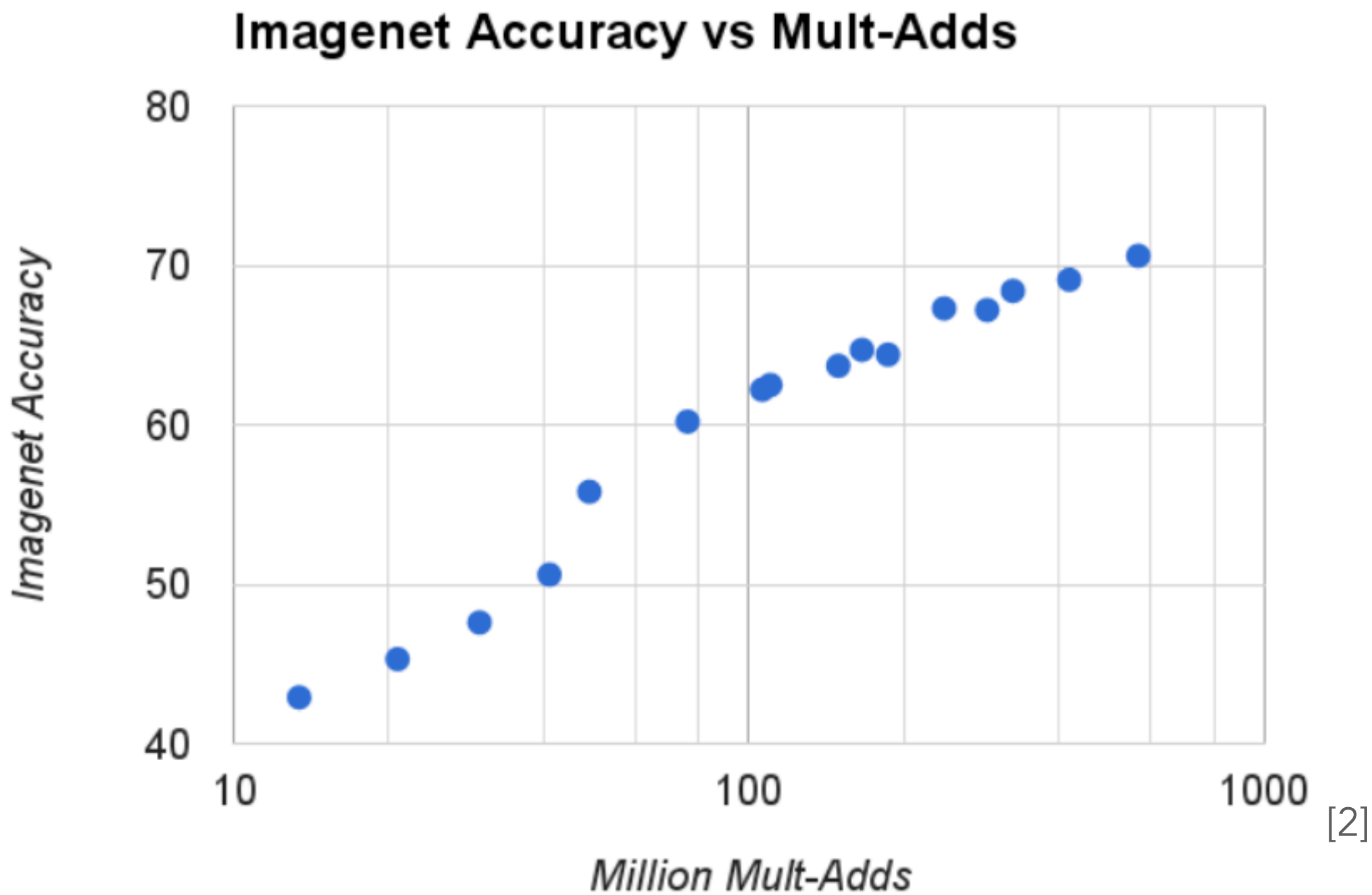
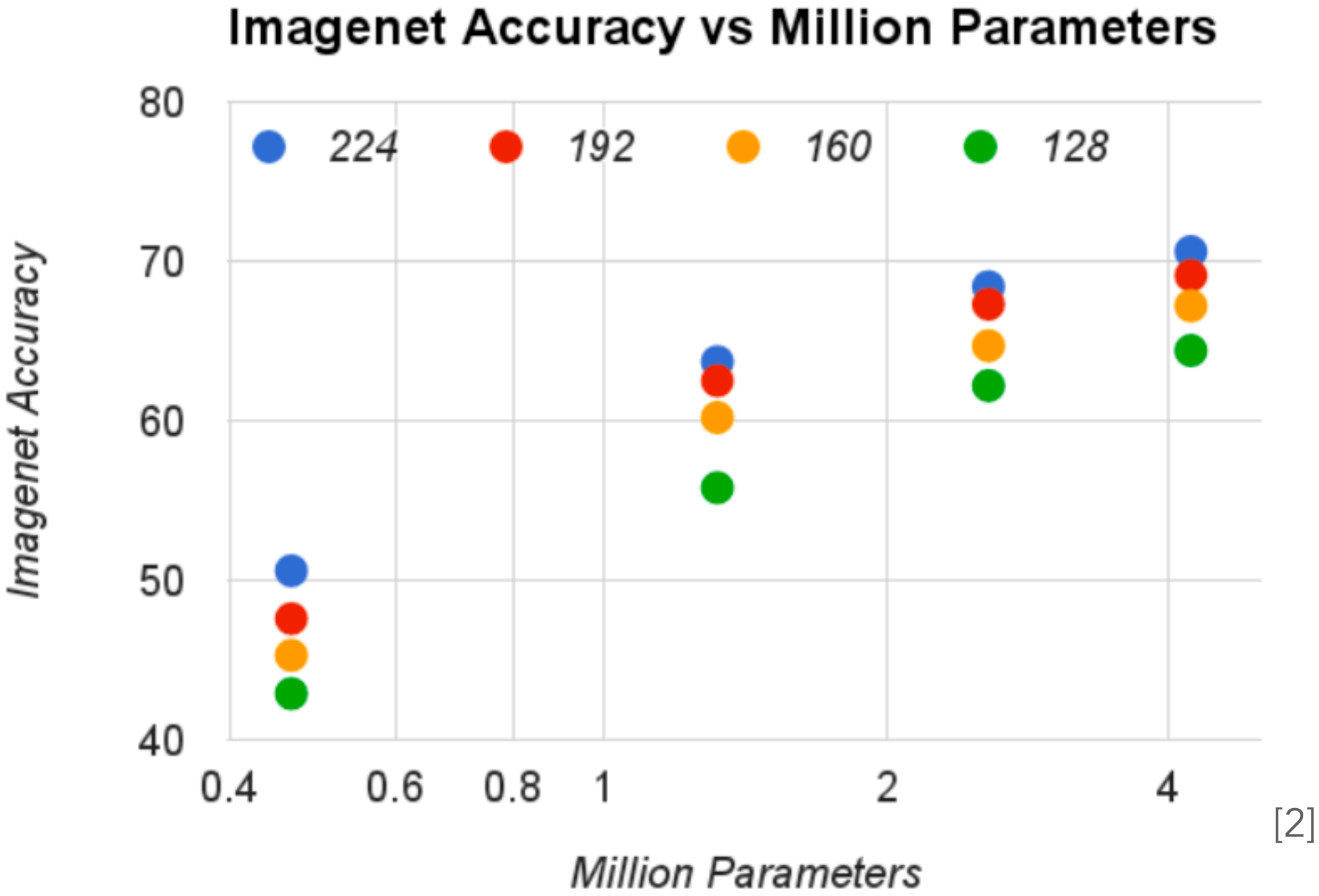
Further Reading:

- A basic tutorial on how to set up an inference backend using python, FastAPI and Tensorflow:
<https://towardsdatascience.com/image-classification-api-with-tensorflow-and-fastapi-fc85dc6d39e8>
- A tutorial on handling image uploads in react via drag and drop:
<https://medium.com/@650egor/simple-drag-and-drop-file-upload-in-react-2cb409d88929>
- A tutorial on using TFJS with react:
<https://levelup.gitconnected.com/build-ad-dog-classifier-with-react-and-tensorflow-js-in-minutes-f08e98608a65>

Additional Slides

How do mobile nets work and how are they helpful?

Model Scaling Parameters

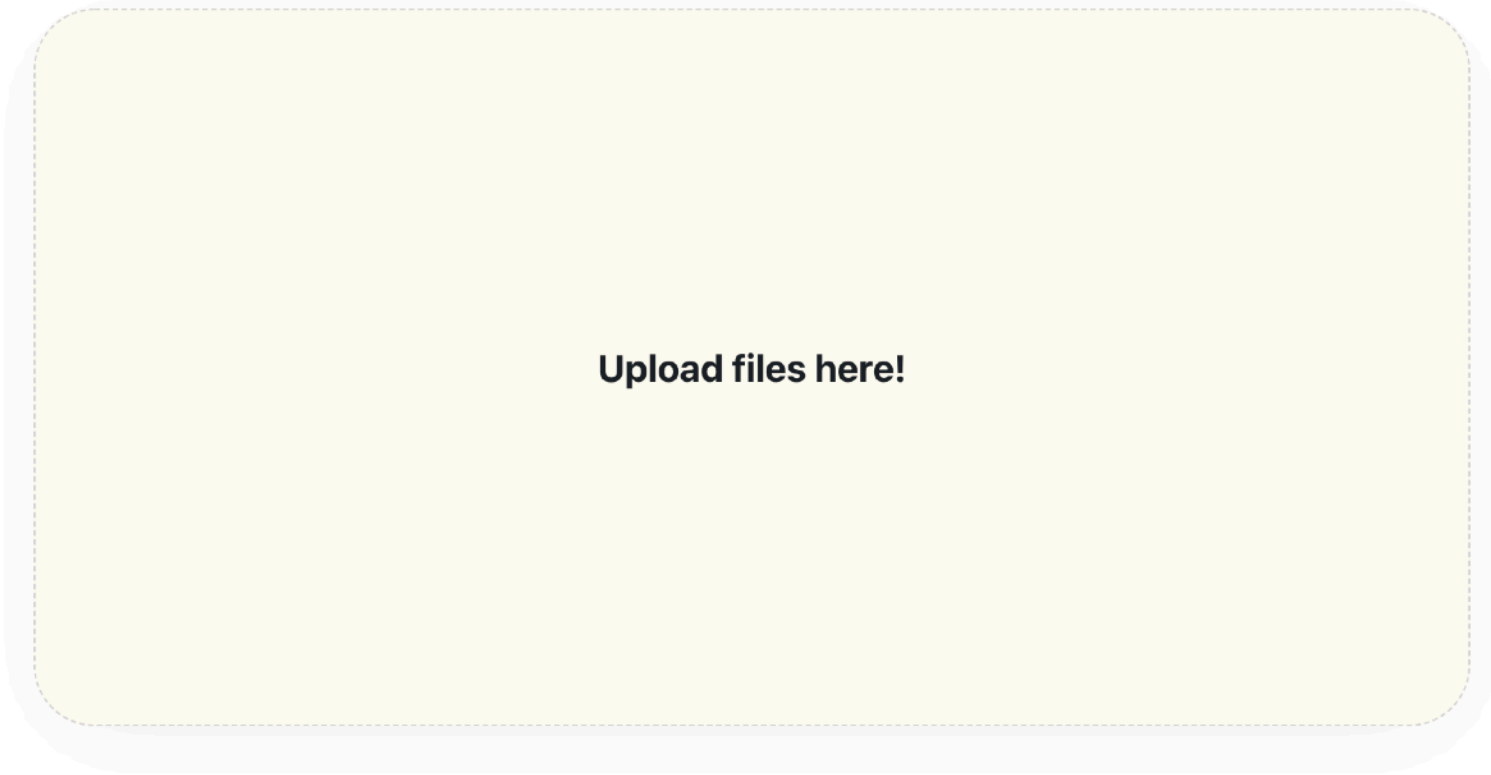


How do I run a model on the browser?

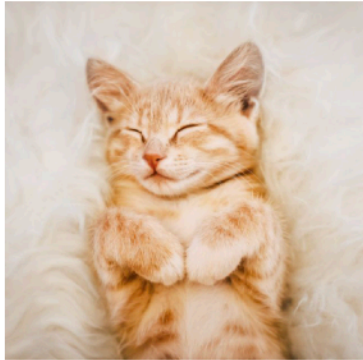
Overview

Server-based Model

Thin Client



llama
58.80 %
hog
24.43 %



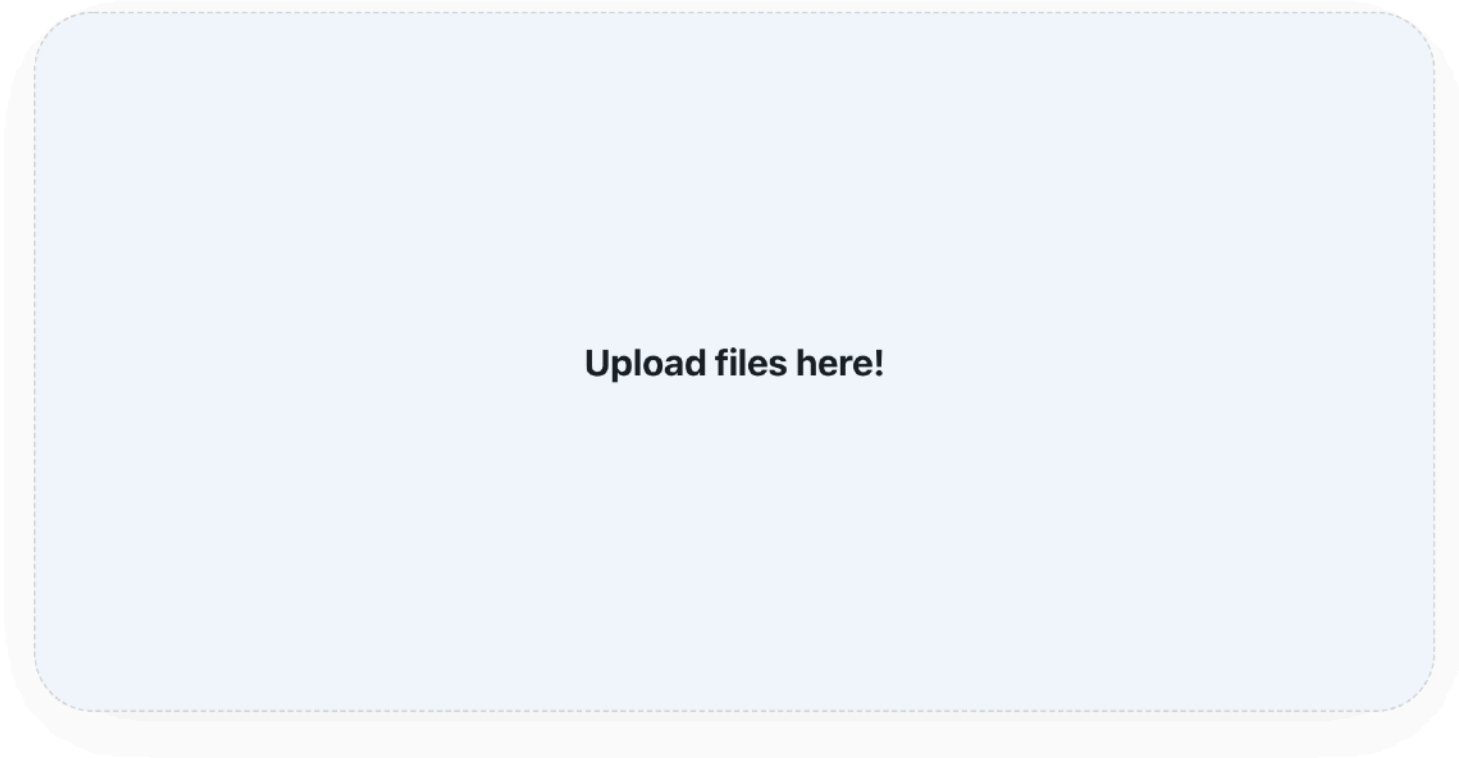
Egyptian_cat
83.68 %
lynx
5.88 %



banana
99.74 %
hook
0.06 %

Client-based Model

Thick Client



llama
0.29
kelpie
0.11
Eskimo dog, husky
0.1



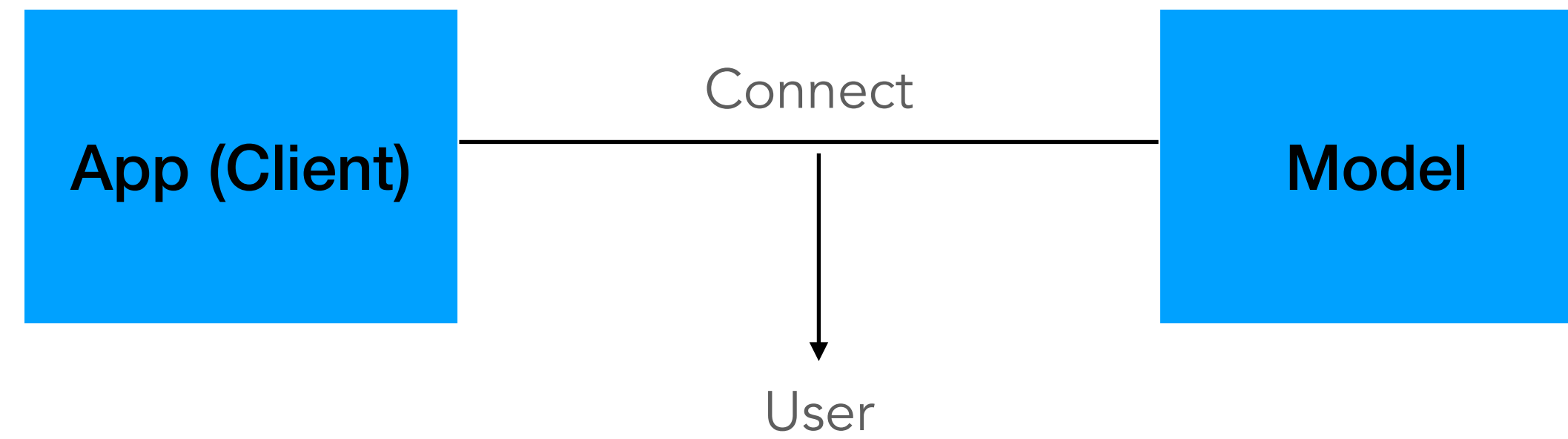
tiger cat
0.33
Persian cat
0.27
tabby, tabby cat
0.14



banana
0.84
nipple
0.14
spaghetti squash
0.01

How do I run a model on the browser?

Overview



1

Write Application
(I/O, Presentation, etc.)

2

Implement Model
(Backend, Parameter
Loading, Inference, etc.)
on the client

3

Make the app available
to the user

How do I run a model on the browser?

Overview

1

Write Application
(I/O, Presentation, etc.)

2

Implement Model
(Parameter Loading,
Inference, etc.)
on the client

3

Make the app available
to the user



How do I run a model on the browser?

Writing a browser application

- ▶ The app has to be able to upload dropped images and display them.
- ▶ There are many viable options to build an application in the browser (for example React, Vue, Angular, native Javascript + HTML + CSS, Typescript, etc...)
- ▶ I used React.js to write the UI.



```
1 import { useState } from "react";
2 import { Upload } from "./upload";
3 import ImageViewer from "./image-viewer";
4
5 const useImgs = (initialState : any[] = [], maxFiles : number = 3) => {
6   const [state, setstate] = useState(initialState);
7   const addImgs = (newDrops) => {
8     const newImgs = newDrops
9       .map((file) => {
10         if (file.type.includes( value: "image")) {
11           file.preview = URL.createObjectURL(file);
12           return file;
13         }
14         return null;
15       })
16       .filter((elem) => elem !== null);
17     setstate([...newImgs, ...state].slice(0, maxFiles));
18   };
19   return [state, addImgs];
20 };
21
22 const ImageClassifier = () => {
23   const [imgs, addImgs] = useImgs( initialState: []);
24   const onFileDrop = (files) => {
25     console.log(files);
26     if (files.length > 0) {
27       addImgs([...files]);
28     }
29   };
30   return (
31     <div>
32       <h2>Thin Client</h2>
33       <Upload onDrop={onFileDrop} />
34       <ImageViewer files={imgs} />
35     </div>
36   );
37 };
38
39 export default ImageClassifier;
```

How do I run a model on the browser?

Integrating a model

- ▶ TFJS(Tensorflow JS)
- ▶ "TensorFlow.js is an open-source hardware-accelerated JavaScript library for training and deploying machine learning models."
<https://github.com/tensorflow/tfjs>
- ▶ To run a mobile net in tfjs requires 2 lines of code



```
import * as mobilenet from "@tensorflow-models/mobilenet";
import * as tf from "@tensorflow/tfjs";
```

```
const App = () => {
  const [model, setModel] = useState( initialState: null);
  useEffect( effect: () => {
    const getModel = async () => {
      const m = await mobilenet.load();
      setModel(m);
    };
    getModel();
  }, deps: []);
```

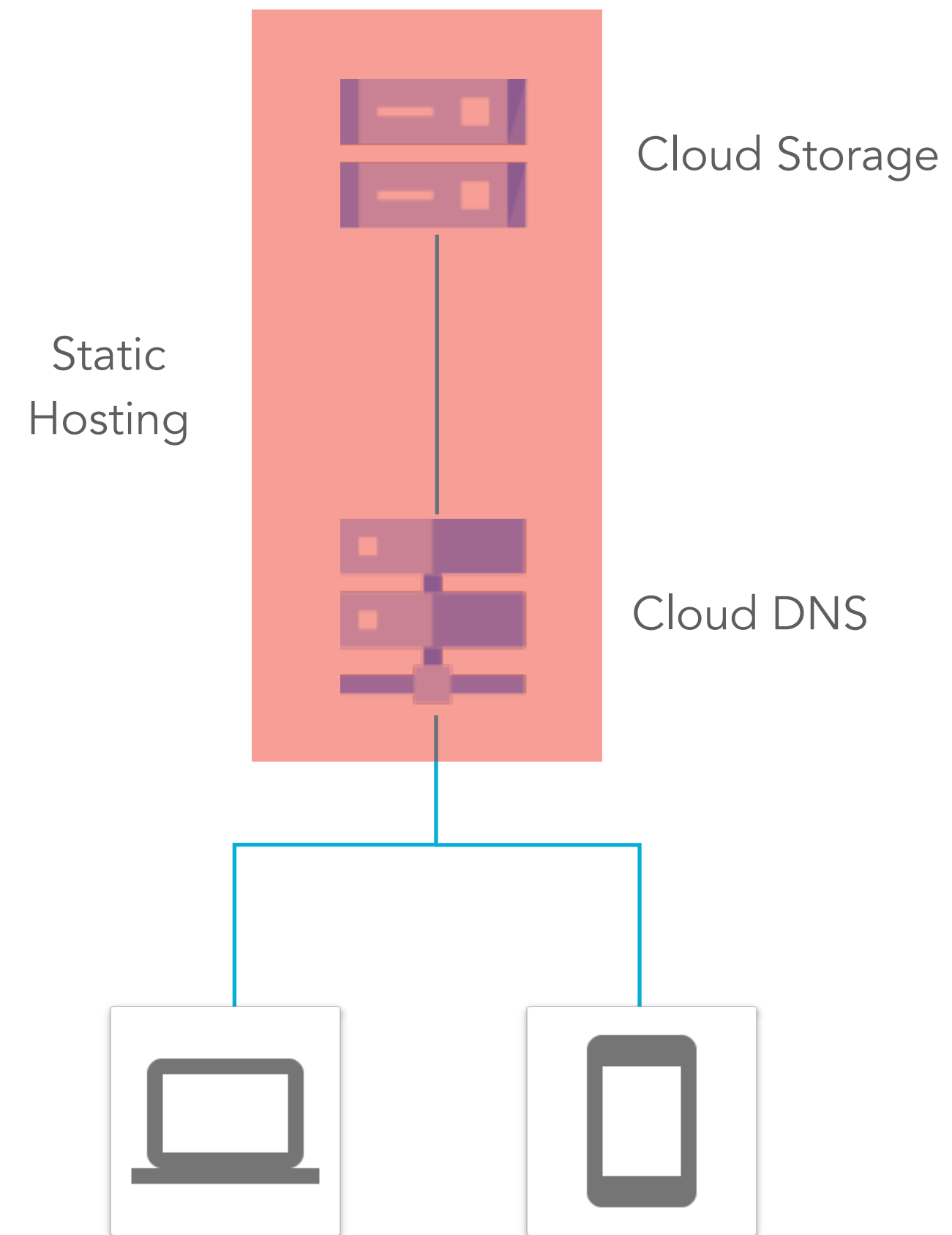
```
const PredictionDisplay = ({ model, img }) => {
  const [predictions, setPrediction] = useState( initialState: null);
  useEffect( effect: () => {
    const fetchClassification = async () => {
      const pred = await model.classify(img);
      setPrediction(pred);
    };
    setPrediction( value: null);
    fetchClassification();
  }, deps: [img]);

  return (
    <PredictionList>
```


How do I run a model on the browser?

Serving a browser CV application

- ▶ Since all dynamic data (i.e. the images to classify) remain at the client, only static data needs to be served.
- ▶ There are lots of Cloud Providers with really affordable (or free) offers for static hosting
- ▶ Examples include Github Pages, Firebase Hosting, ...



How do I run a model on the browser?

What if I want to use a thin client?

```
1 import axios from "axios"
2
3 const request = axios.create({
4   baseURL: 'http://localhost:8000',
5   timeout: 10000,
6 });
7
8 const getPrediction = async (image) => {
9   let blob = await fetch(image.preview).then(r => r.blob());
10  const formData = new FormData();
11  formData.append('file', blob);
12  return request.post('/predict/image', formData,
13    { config: { headers: { "Content-Type": "multipart/form-data" } } })
14 }
15
16 export {getPrediction}
```

We need to add a request middleware to the web app ...

```
1 import uvicorn
2 from fastapi import FastAPI, File, UploadFile
3 from serve_model import predict, read_imagefile
4
5
6 app = FastAPI()
7 @app.post("/predict/image")
8 async def predict_api(file: UploadFile = File(...)):
9   print(file.filename)
10  image = read_imagefile(await file.read())
11  prediction = predict(image)
12  return prediction
13
14
15 if __name__ == "__main__":
16   uvicorn.run(app, debug=True)
```

...,which calls the API endpoint of a python-based server.

