

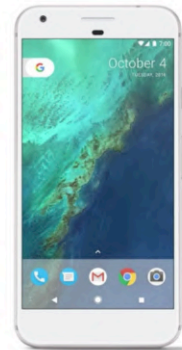
Federated Learning

Data is born at the edge

Billions of phones & IoT devices constantly generate data

Data enables better products and smarter models

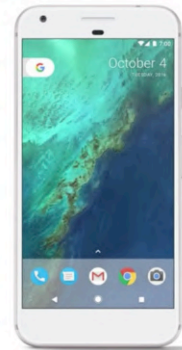
With everyone becoming more privacy aware, we want to be nice if you can train a model without needing to upload highly identifiable personal information to a Cloud hosting server. What if instead, data can be covered anonymous or some of the training can even be done in a distributive way.



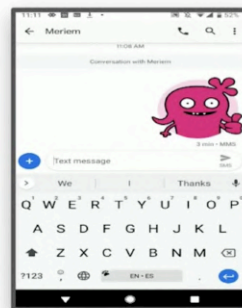
Can data live at the edge?

On-device inference offers:

- Improved latency
- Works offline
- Better battery life
- Privacy advantages



Gboard: mobile keyboard





Gboard machine learning

Models are es

- tap typing
- gesture typii
- auto-correct
- predictions
- voice to text

Gboard machine learning

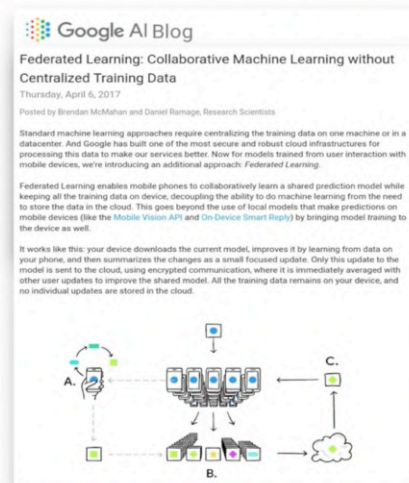
Models are essential for:

- tap typing
- gesture typing
- auto-corrections
- predictions
- voice to text
- and more...

All on-device for latency and reliability



<https://ai.googleblog.com/2017/04/federated-learning-collaborative.html>



Not Available



Available



Not Available

Examples of the learning and not the data used to generate the learning can be shared across all users. So how does all of this work?

Well, let's consider a simple scenario like this one. I have a population of users each with a mobile device and each with my model running on it. If I want to do some training on these devices, I realize that I'm getting them to do something that's computationally expensive.

So I can select a subset of them that are available in the sense that they're not currently being used, they're plugged in, and they're charging and they're idle. The principle here is that I don't want to damage the user experience. Of the available ones, some of them will have data that's relevant to the problem I'm trying to solve like training to understand the keyboard usage better for example. So you'll pick a subset of the available devices.



Not Available



Available

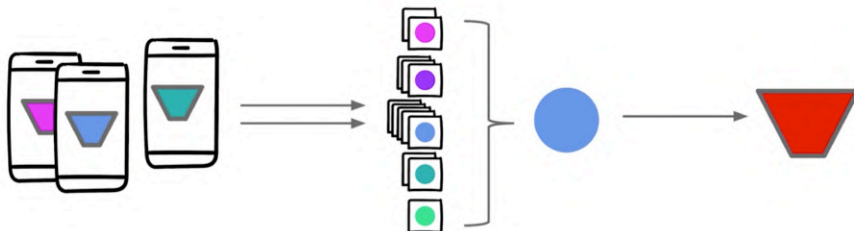


Not Available

These will then receive a training model which can then be retrained on the device



The results of the training, not the data used to perform the training, is sent to the server.



The server then uses this to retrain the master model.

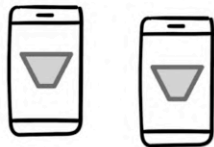
Before deploying the new model to your customers devices, of course, the model should be tested and what better for testing than to use a very similar approach using a population of the user set to test if the model works well with their data.



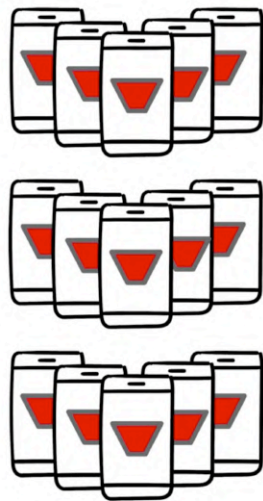
Available



So we can take a subset for testing the same way as we took a subset for training and deploy a new model to those.



Using a process like this and iterating, then over time we can improve the model significantly and the user's personal data never leaves their device. Each device can benefit from the experience of others so the learning is distributed across all users and the term federated is used to describe this.



Previously, you had an overview of the concepts of federated learning and how a master model could be retrained with distributed users data without the user's data being uploaded to a server where it could potentially be misused. The data always stays on the user's device. At the end of that demonstration, we realized that there is the potential for the retrained models that the user has uploaded to possibly be reverse engineered to get at that data.

So with that in mind, I'd like to go through two concepts now that can help ease any fears here. The first is that by principle federated learning works only on aggregates and the second to show how the data can be encrypted by the clients on route to the server. First, let's take a look at an aggregation methodology that can be used to maintain an individual user's privacy.

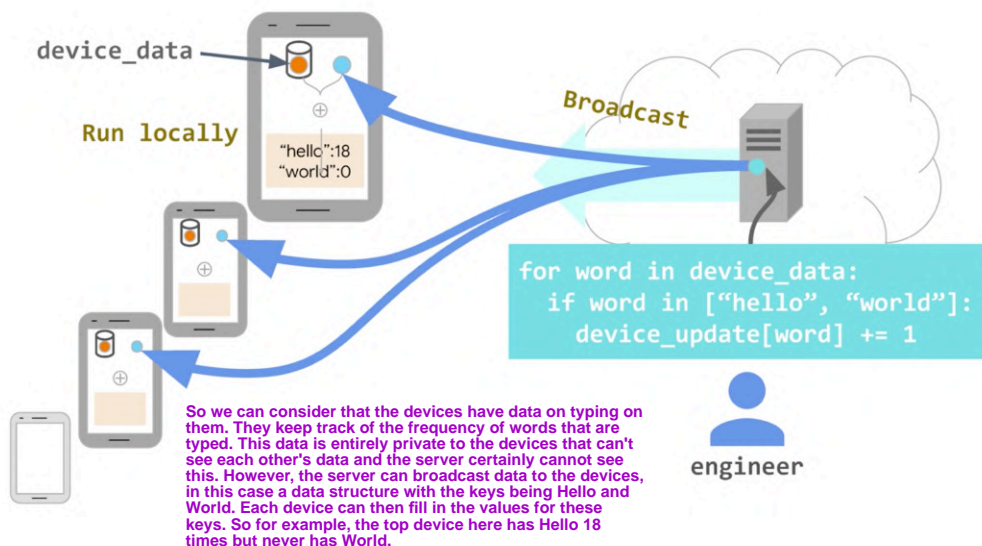
Hello World
Andrew
Mom
Who's there?
Hello I'm in a place called Vertigo
...

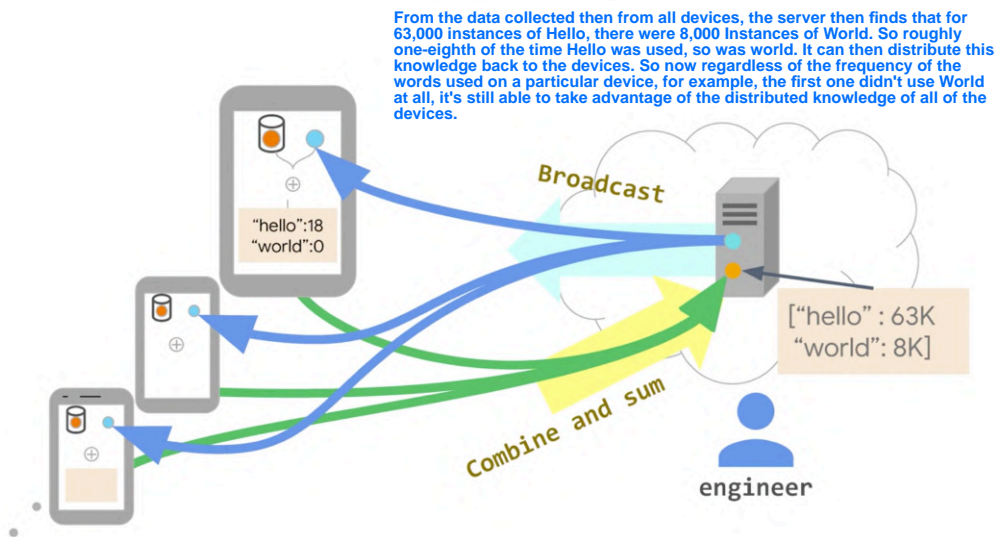
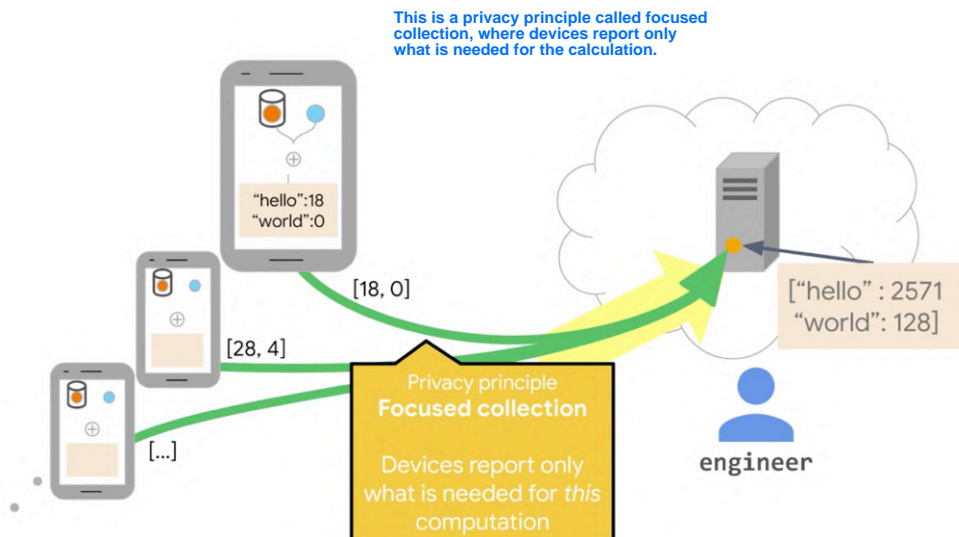
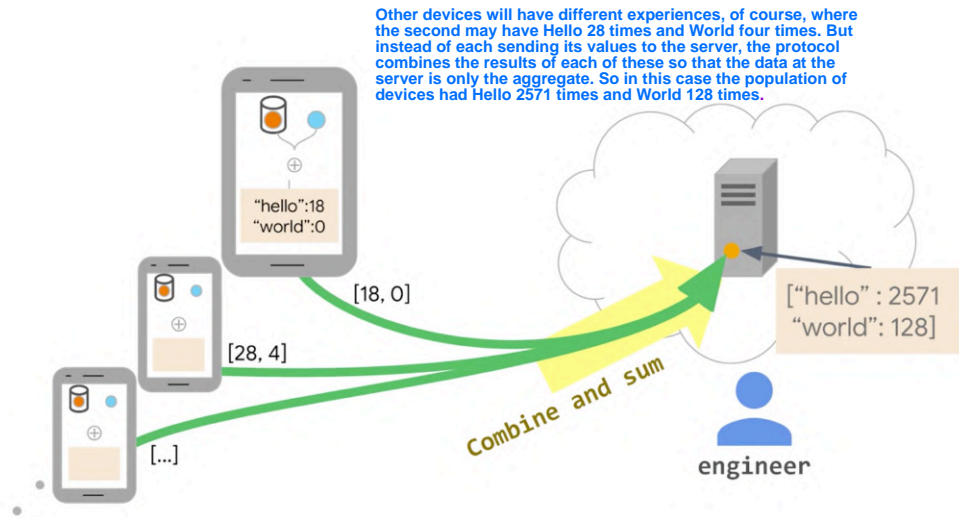
Consider this scenario we're retraining a keyboard predictor and right now it doesn't know what word to use to follow Hello, but our users generally follow the word Hello with a word World. Hello world. So we'd like a result like this where if they type Hello then maybe the most common words that follow hello would be listed.

Hello World

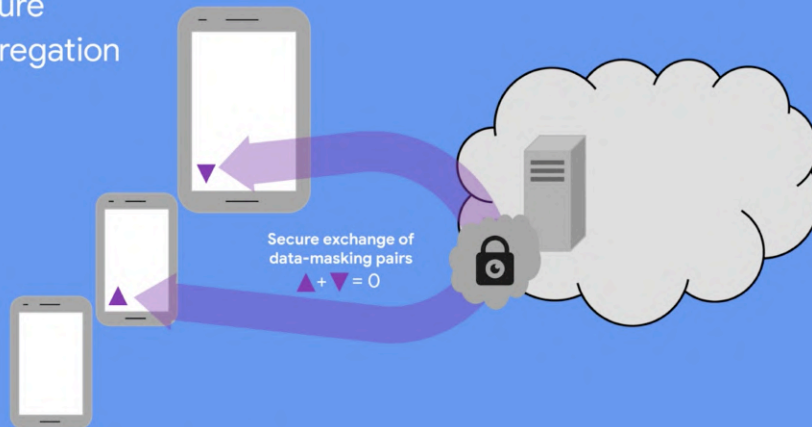
Follows 'Hello'
>10% of the time

But let's keep it really simple and just theorize how often World will follow Hello as people type. We think it happens a lot. But let's learn from our users just how often it follows.



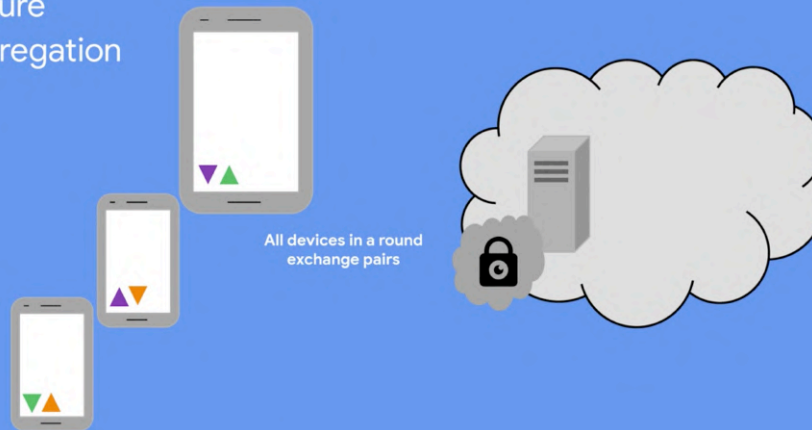


Secure aggregation



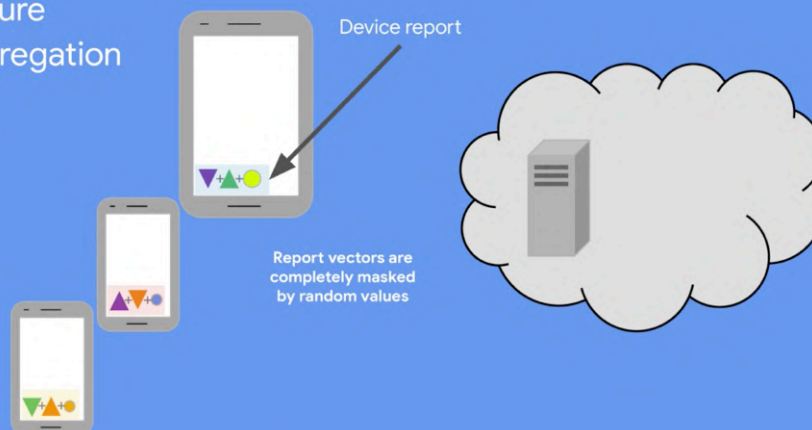
K. Bonawitz, et al. Practical Secure Aggregation for Privacy-Preserving Machine Learning. CCS 2017.

Secure aggregation



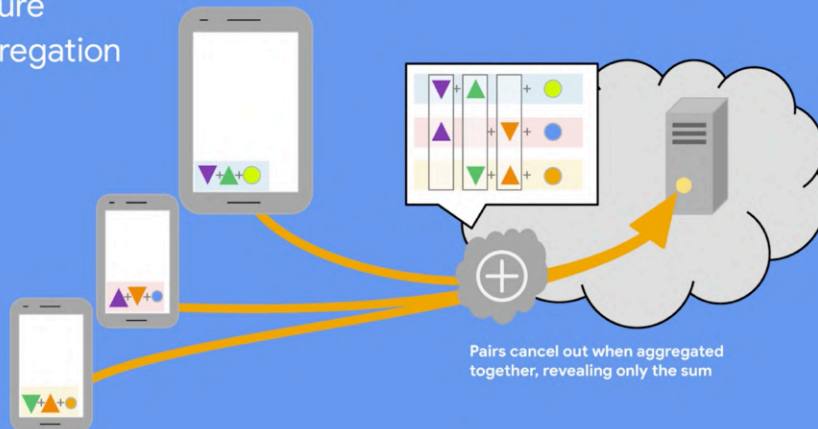
K. Bonawitz, et al. Practical Secure Aggregation for Privacy-Preserving Machine Learning. CCS 2017.

Secure aggregation



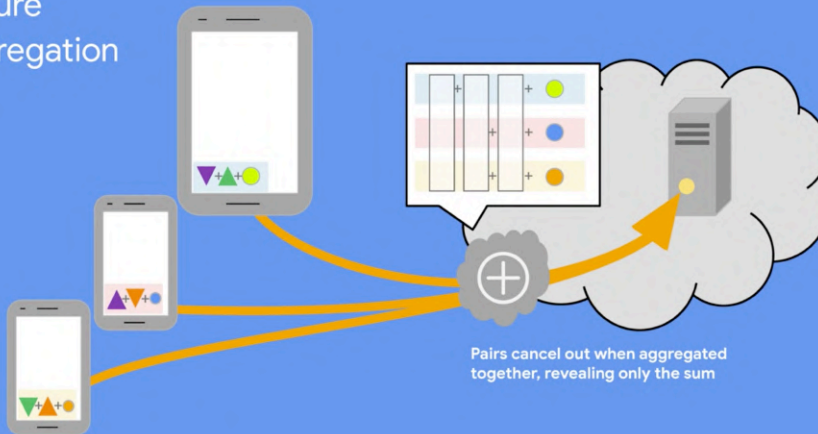
K. Bonawitz, et al. Practical Secure Aggregation for Privacy-Preserving Machine Learning. CCS 2017.

Secure aggregation



K. Bonawitz, et al. Practical Secure Aggregation for Privacy-Preserving Machine Learning. CCS 2017.

Secure aggregation



K. Bonawitz, et al. Practical Secure Aggregation for Privacy-Preserving Machine Learning. CCS 2017.

<https://eprint.iacr.org/2017/281.pdf>

Practical Secure Aggregation for Privacy-Preserving Machine Learning

Kath Bonawitz¹, Vladimir Iyengar², Ben Kreuter²,
Antonio Mardelloski¹, J. Brandon McDaniel¹, Suresh Patel¹,
Daniel Ramage¹, Aaron Segal¹, and Kara Seth¹
¹Google, Mountain View, CA 94043
²Google, Mountain View, CA 94043
kardones@cs.cornell.edu

Cornell Tech, 2 West Loop Rd., New York, NY 10044

1 INTRODUCTION

Machine learning models trained on sensitive real-world data are increasingly being used to make decisions that affect people's lives. In the medical domain, for example, machine learning models are used to predict disease outcomes [18]. And the widespread use of mobile devices means that there is a vast amount of sensitive data in the hands of users.

However, large-scale collection of sensitive data entails risks. A particularly high-profile example of the consequences of collecting sensitive data occurred in 2016, when the video history of a woman for the US Supreme Court was published without her consent [15]. The fact that the data was collected in the first place is a concern, but the fact that it was published without her consent is a more serious one. It is a reminder that sensitive data is becoming available [27].

There are two main approaches to addressing privacy-preserving machine learning for sensitive data: federated learning [19] and secure aggregation [28]. Federated learning involves training a model on multiple devices, and then aggregating the results. Secure aggregation involves training a model on a central server, and then aggregating the results. In this paper, we focus on secure aggregation.

The problem of aggregating a collection of data in a secure manner is a well-studied problem in cryptography. The problem of aggregating a collection of data in a secure manner is a well-studied problem in cryptography. The problem of aggregating a collection of data in a secure manner is a well-studied problem in cryptography.

We are particularly focused on the setting of mobile devices, where communication is inherently expensive, and devices are resource-constrained. Given these constraints, we would like our protocol to have as few communication rounds as possible, and to be as simple as possible to implement. In this paper, we describe a secure aggregation protocol that meets these goals.

1.1 Our Results

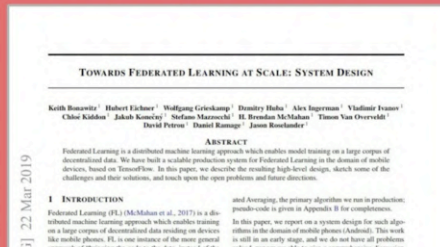
We present a protocol for securely aggregating sums of vectors, which has a constant number of rounds, low communication overhead, robustness to failures, and which requires only one server with trusted setup. In our design, the server has no access to the data being aggregated, and it computes the final result. We present two variants of the protocol: one for aggregation of sums of vectors, and one for aggregation of sums of scalars. Both variants are secure against active adversaries, and both are secure against passive adversaries.

1.2 Organization

In Section 2 we describe the machine learning application that motivates this work. In Section 3 we review the cryptographic primitives we use in our protocol. We then proceed to give a high-level overview of our protocol design in Section 4, followed by a formal protocol description in Section 5. In Section 6 we prove security against honest-but-curious (passive) adversaries and include a high-level discussion of privacy against active adversaries. In Section 7 we give performance results for our protocol, and in Section 8 we conclude with a discussion of related work in Section 9.

Google's federated system

Towards Federated Learning at Scale: System Design
ai.google/research/pubs/pub47976



K. Bonawitz, et al. Towards
Federated Learning at Scale:
System Design. SysML 2019.

tensorflow.org/federated

TensorFlow Federated

What's in the box

Federated Learning (FL) API

- Implementations of federated training/evaluation
- Can be applied to existing TF models/data

Federated Core (FC) API

- Allows for expressing new federated algorithms

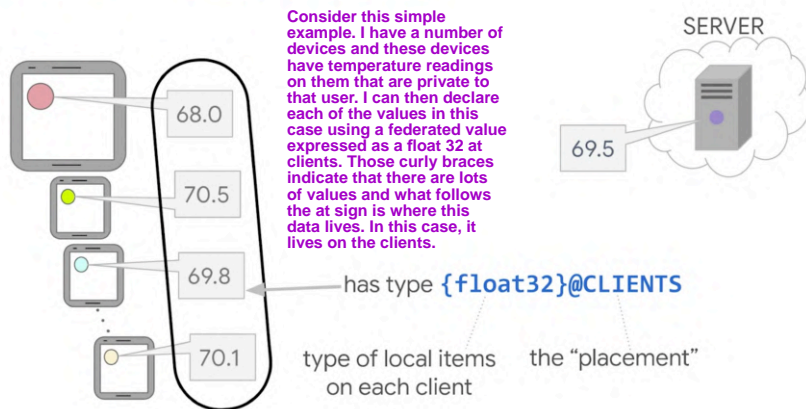
Local runtime for simulations

TensorFlow Federated offers two main APIs for federated learning. The federated learning API contains implementations of federated training and evaluation that can be applied to existing core as models. So you can experiment with learning using them.

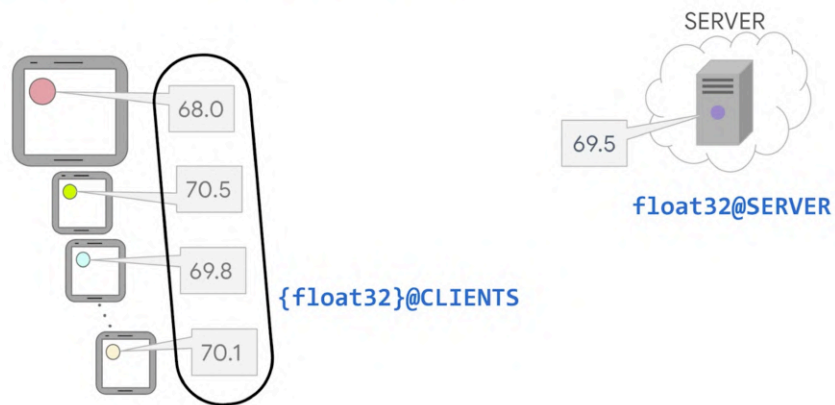
Note that there are no more mobile APIs yet. Everything is done in simulation and the collabs at the TensorFlow site will demonstrate this.

Similarly, the federated core API allows you to express new federated algorithms and to test them in a simulated environment.

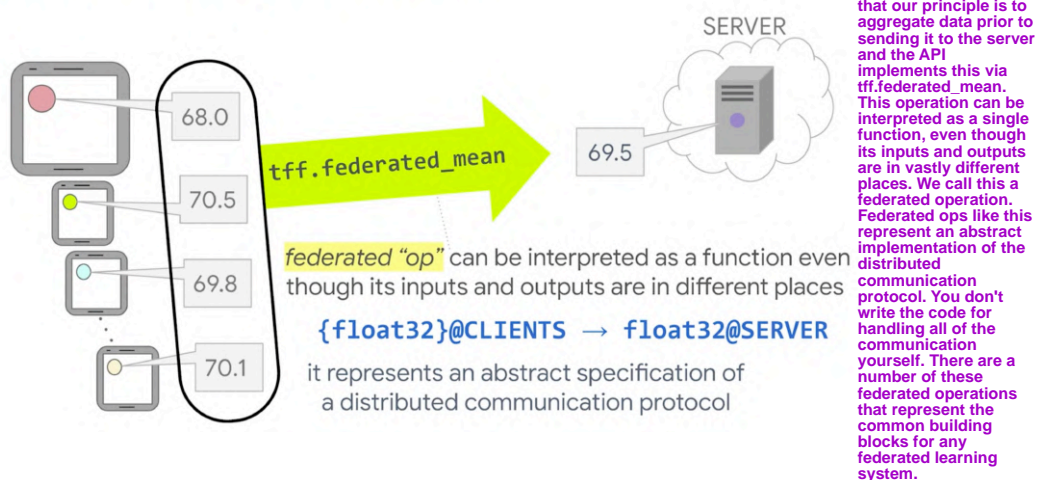
Federated computation in TFF



Federated computation in TFF



Federated computation in TFF



Federated computations in TFF

```
READINGS_TYPE = tff.FederatedType(tf.float32, tff.CLIENTS)
```

```
# An abstract specification of a simple distributed system
@tff.federated_computation(READINGS_TYPE)
def get_average_temperature(sensor_readings):
    return tff.federated_mean(sensor_readings)
```

Let's look at a brief code example using TFF. I'm not going to go to into depth so it might look a little confusing, but the TensorFlow site has a colab that you can walk through that implements all of this. We'll first implement a federated type that represents inputs. Note that we specify that it's on our clients.

Federated computations in TFF

```
READINGS_TYPE = tff.FederatedType(tf.float32, tff.CLIENTS)
```

```
# An abstract specification of a simple distributed system
@tff.federated_computation(READINGS_TYPE)
def get_average_temperature(sensor_readings):
    return tff.federated_mean(sensor_readings)
```

We're then going to implement a function to get the average of the sensor readings. But in order for this to be federated, we'll use a function decorator saying that it's a federated computation using `@tff.federated_computation`

Federated computations in TFF

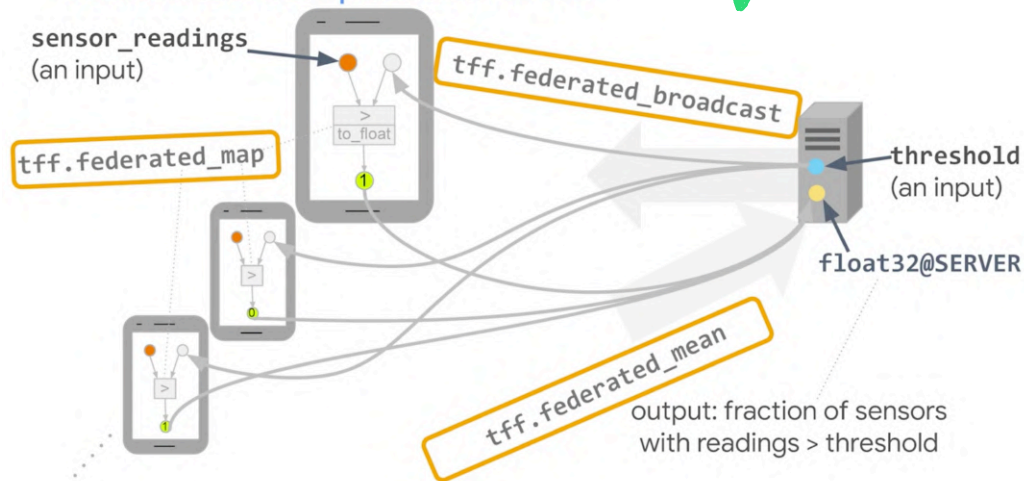
```
READINGS_TYPE = tff.FederatedType(tf.float32, tff.CLIENTS)
```

```
# An abstract specification of a simple distributed system
@tff.federated_computation(READINGS_TYPE)
def get_average_temperature(sensor_readings):
    return tff.federated_mean(sensor_readings)
```

Similarly, if we want to do a federated computation, we will broadcast from the server a value that we want to learn about. For example, if the sensor readings might be above a certain value. We can do this with `tff.federated_broadcasts`. In this case, we want to find the number of machines where the value of their temperature is above a certain amount, say 70 degrees. The server will broadcast the data and the device will check it against its internal value giving us a one if it's higher and a zero if it isn't. These values are then aggregated using a federated mean and the server can find out how many of the devices are above that value overall without it ever knowing

In this function I'll return `tff.federated_mean` of the sensor readings. Under the hood, the readings are being gathered by the server using the secure protocols, decrypted, and then returns to me as an aggregates.

Federated computation in TFF



```
THRESHOLD_TYPE = tff.FederatedType(
    tf.float32, tff.SERVER, all_equal=True)

@tff.federated_computation(READINGS_TYPE, THRESHOLD_TYPE)
def get_fraction_over_threshold(readings, threshold):

    @tff.tf_computation(tf.float32, tf.float32)
    def _is_over_as_float(val, threshold):
        return tf.to_float(val > threshold)

    return tff.federated_average(
        tff.federated_map(_is_over_as_float, [
            readings, tff.federated_broadcast(threshold)]))
```

The code for this is very similar. We'll create a federated type for a threshold and we'll indicate that it lives on the server with `tff.server`. We'll then decorate a function called `get_fraction_over_threshold` as a federated computation.

```
THRESHOLD_TYPE = tff.FederatedType(
    tf.float32, tff.SERVER, all_equal=True)

@tff.federated_computation(READINGS_TYPE, THRESHOLD_TYPE)
def get_fraction_over_threshold(readings, threshold):

    @tff.tf_computation(tf.float32, tf.float32)
    def _is_over_as_float(val, threshold):
        return tf.to_float(val > threshold)

    return tff.federated_mean(
        tff.federated_map(_is_over_as_float, [
            readings, tff.federated_broadcast(threshold)]))
```

This function will return a federated mean which takes the readings and the thresholds that was broadcast as inputs to a federated map value

```

THRESHOLD_TYPE = tff.FederatedType(
    tf.float32, tff.SERVER, all_equal=True)

@tff.federated_computation(READINGS_TYPE, THRESHOLD_TYPE)
def get_fraction_over_threshold(readings, threshold):

    @tff.tf_computation(tf.float32, tf.float32)
    def _is_over_as_float(val, threshold):
        return tf.to_float(val > threshold)

    return tff.federated_average(
        tff.federated_map(_is_over_as_float, [
            readings, tff.federated_broadcast(threshold)]))

```

The computation is to simply return if the val of the client is over the threshold value.

TensorFlow Federated: Machine Learning on Decentralized Data

TensorFlow Federated (TFF) is an open-source framework for machine learning and other computations on decentralized data. TFF has been developed to facilitate open research and experimentation with Federated Learning (FL), an approach to machine learning where a shared global model is trained across many participating clients that keep their training data locally. For example, FL has been used to train prediction models for mobile keyboards without uploading sensitive typing data to servers.

TFF enables developers to simulate the included federated learning algorithms on their models and data, as well as to experiment with novel algorithms. The building blocks provided by TFF can also be used to implement non-learning computations, such as aggregated analytics over decentralized data. TFF's interfaces are organized in two layers:

- Federated Learning (FL) API**
This layer offers a set of high-level interfaces that allow developers to apply the included implementations of federated training and evaluation to their existing TensorFlow models.
- Federated Core (FC) API**
At the core of the system is a set of lower-level interfaces for concisely expressing novel federated algorithms by combining TensorFlow with distributed communication operators within a strongly-typed functional programming environment. This layer also serves as the foundation upon which we've built Federated Learning.

TFF enables developers to declaratively express federated computations, so they could be deployed to diverse runtime environments. Included with TFF is a single-machine simulation runtime for experiments.

```

import tensorflow as tf
tf.compat.v1.enable_v2_behavior()
import tensorflow_federated as tff

# Load simulation data.
source, _ = tff.simulation.datasets.emnist.load_data()
def client_data(n):
    return source.create_tf_dataset_for_client(source.client_ids[n]).map(
        lambda e: {
            'x': tf.reshape(e['pixels'], [-1]),
            'y': e['label'],
        }).repeat(10).batch(20)

# Pick a subset of client devices to participate in training.
train_data = [client_data(n) for n in range(3)]

# Grab a single batch of data so that TFF knows what data looks like.
sample_batch = tf.nest.map_structure(
    lambda x: x.numpy(), iter(train_data[0]).next())

# Wrap a Keras model for use with TFF.
def model_fn():
    model = tf.keras.models.Sequential([
        tf.keras.layers.Dense(10, tf.nn.softmax, input_shape=(784,),
            kernel_initializer='zeros')

```

<https://www.youtube.com/watch?v=89BGjQYA0uE>



Intermediate

Federated Learning

