2024.12.17

**design principle:**
**progressive disclosure of complexity**

**only one per import**

**works with NumPy, Pandas, Tensorflow** `Dataset`**, PyTorch** `DataLoader` **regardless of backend**

**models are similar to layers but can have component models**

**Backends: JAX, Tensorflow, PyTorch**

**packaged with popular dataset** ~~**downloaders**~~ *pipelines*

**the layer abstraction**

**Keras first impressions**

**model-centric**

`init` **with configuration,** `build` **if it has persistent data/params depending on input shape,** `call`

**with** `Input` **layer builds automatically, without it build manually with batch input shape**

```
model = keras.Sequential([   ...list of layers...   ])
model.compile(loss, optimizer, metrics)
callbacks =
   [ModelCheckpoint(..), EarlyStopping(..)]
model.fit(x_train, y_train, batch_size,
               epochs, validation_split, callbacks)
score = model.evaluate(x_test, y_test)
model.save("final_model.keras")
predictions = model.predict(x_test)
```

2024.12.18

**reusing a layer expression in different models does not share weights, in same model shares weights**

**in OCANNL, a tensor expr. function shares weights, a layer / block with `~config` does not**

**composing with a model shares weights**

<span style="color:blue">upcoming:
- training and eval
- distributed training</span>

**model inputs can be a list, outputs can be a dictionary**

**functional API = layer expressions instead of `Sequential`**

**saved model includes:**
**- architecture (layer expression)**
**- weight values (params)**
**- training config**
**- optimizer and its state**

**for cyclic or recursive computations: subclass `Model`**

**Keras styles: Sequential, functional, subclassing**

**can mix-and-match Sequential, layer expressions and subclassing -- via composing (sub)models**

**auto-propagated `call` args**

**`mask`**

**`training`**

**layers and models have a `trainable` flag**

**regenerated per-call**

**bool tensor if model input shape**

**train vs. inference, handled by built-in train, eval, predict loops**

**individual weights can also be non-trainable**

**layers can `add_loss` to models that use them**

2024.12.19

sample weights: per-sample influence on loss

class weights: balance classes without resampling

Data sources i.e. input pipelines are iterator-based (except NumPy), offer batching and shuffling, keras-specific one is mulitcore.

Dynamic learning rate schedules are callbacks that modify the optimizer.

Callback class has methods specific to: begin/end of whole/batch/epoch of train / test i.e. eval / predict i.e. infer

Keras training

For saving/loading, custom layers etc. must define get_config, usually captures init arguments.

Ideas for callbacks: checkpointing, early stopping, changing learning rate when plateau, fine-tuning of top layers when plateau, emailing on performance thresholds, TensorBoard, CSVLogger.

allows e.g. subclassing a GAN model

Progressive intervention into a model's training:
- override train_step and/or test_step (of eval) using model's forward-call and loss initerface;
- as above but inline loss;
- write the training and/or eval loop from scratch.

examples generate the derivative at each train step

JAX example jit-compiles the full train step

2024.12.20

**OCANNL's** `DeviceMesh : dev`

In OCANNL, better fit to link DeviceMesh with a routine rather than a tensor.

grid configured manually but sharding done by program search

**same as** `tensorflow.dtensor`

per-cluster mesh config passed to the mesh backend functor

`DeviceMesh`

`TensorLayout`

organizes devices into N-dim grid with axis_names

**Keras distributed**

assigns axes of any tensor (positionally) to sharded on a given mesh axis, or replicated

is synchronous

no events

tied to a device_mesh (might initially be unset)

`DataParallel`

`ModelParallel`

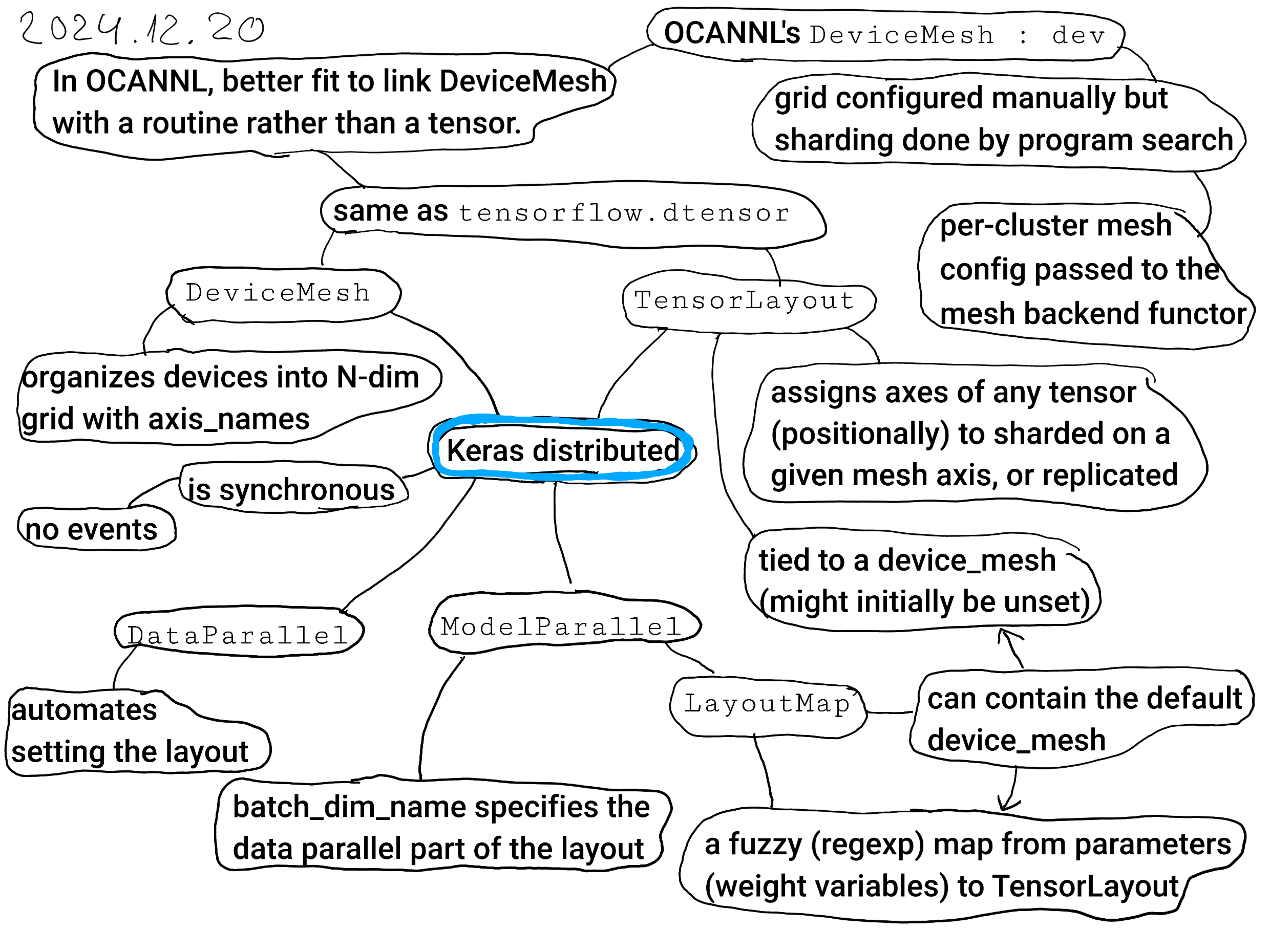`LayoutMap`

can contain the default device_mesh

automates setting the layout

batch_dim_name specifies the data parallel part of the layout

a fuzzy (regexp) map from parameters (weight variables) to TensorLayout

counter-based PRNGs are better for parallelism

unassigned input axes are replicated / tiled as in DTensor

partitions tensors preserving the ranks (i.e. nums of axes)

sharding mesh + PartitionSpec (like TensorLayout) = device-like

`Array`: like DTensor

JAX distributed

`shard_map` takes a mesh and partition specs for inputs and output

inferred layout of outputs minimizes copying

layout propagation / inference

mapped func result shape must have rank sufficient for concatenation of shareding axes in output partition spec

error when explicit layouts of inputs disagree

default layout inputs can be moved and resharded automatically to fit other inputs

unassigned output axes are un-replicated: result is selected from just a subset of devices, assuming that it's the same on other groups of devices

`with_sharding_constraint` redirects layout inference

caller can pick mesh axes that are propagated rather than set manually on inputs / output

2024.12.22

**JAX collectives**

transposes blocks along an on-device and a cross-device axis
`all_to_all`

concatenates blocks along an axis, replicating a tensor
`all_gather`

replicates the summed axis
`psum`

sends tensor(s) by permuting a mesh axis
`ppermute`

= ppermute + add, no replication

communicate across devices from within `shard_map`

`psum_scatter`

to overlap comp. and comm. reshape to add an axis and loop over it inside the map

for best shift perf on TPUs, split blocks in half and shift bidirectionally

if not overlapped by XLA

**NN parallel patterns in JAX**

SPMD pipeline parallel

data parallel

FSDP

tensor parallel

shard data and params on corresponding features axis, `psum_scatter` activations

for same structure layers: `shard_map` over concatenated params, `ppermute` to advance the pipeline

shard data on a batch axis, `pmean` the loss

also shard params, on the batch mesh axis

all_gather inside predict, jax.remat to re-gather on backward pass

FSDP + TP

other sharding is automatic

explicit psum for features (in TP automatic sum->psum)

2024.12.23

processes must agree on per-device sizes

control flow must not diverge, watch out:
length of training loop, iteration order

very restrictive approach:
- SPMD: all processes same computations
- all processes same number of devices
- all devices the same (e.g. H100)

death of any process kills others

but allows running `shard_map`
etc. without changes

each JAX process runs independently,
no one controller but one coordinator

**JAX distributed multi-host**

NVIDIA backend: Collective
Communications Library NCCL

sometimes the storage locality
disagrees with computation
locality -- load `jax.Array` with
storage sharding, and add
`with_sharding_constraint`
for efficient resharding

JAX integrates with tf.data.Dataset

2024.12.26

**XLA instruction set**

`BatchNormTraining/Grad/Infer`

`ConvWithGeneralPadding`

`Scatter`, `SelectAndScatter`: **non-deterministic loop of updates**

`Conditional`

`While`

**domain- or algo-specific**

`Fft` **forward and inverse Fourier**

`OptimizationBarrier`

**control-flow-like**

`Cholesky`, `TriangularSolve`

`AfterAll` **for sequencing (like tensor-centric events)**

`Clamp` **to min/max**

`XlaOp` = **tensor**

`CompositeCall`: **to define composite functions**

**cross-replica:** `AllGather`, `AllReduce`, `AllToAll`, `CollectivePermute`, `ReduceScatter`

**can define asynchronouns funcs: start, update loop, done**

`Infeed`: **reads a tensor from an implicit channel on a device**

**persisted autotuning: cache on disk for speed and determinism**

**vectorized:** `Reduce`, `Map`

`Iota`: **constant literal initialized on device without transfer**

`Recv` **and** `Send`: **communicate via shared channel**

**tensor structure**

`Transpose`: **permute axes**

`Gather` **general idea: convert a list of offsets into tensors into a tensor with a new batch dimension**

`Collapse`

**Also arithmetic**

`Broadcast`

`Concatenate`

2024.12.27

3 compilation routes: libraries like cuBLAS & cuDNN; tiling followed by Triton; Emitters

Transpose and Reduction emitters, using shared memory

two loops: coalesced reads to shared mem; then `sync_threads`; then coalesced writes

Partitioning: tensors are emitted in a single function when they interact pointwise without duplication.

Loop emitter is default (no "hero")

XLA Emitters

Subkernel function inputs: "inflow" tensors and indices of "outflow" tensors; outputs: "outflow" values at the indices. Kernel function: takes both "inflow" and "outflow" tensor args.

Other emitters: Concatenate, Dynamic Update Slice, Input slices, Scatter

symbolically computes indexing maps between tensors, e.g. input<->output

Only single-call functions are inlined.

for reasoning on mem. coalescing and tiling propagation

loop traversals linear in output tensors for coalesced writes, with boundary checks inside

for emitting index transformations (transpose, broadcast, reshape, slice, reverse)

tensors flattened to 1D as in memory

loop unrolling

only contiguous accesses get inlined as transfer reads