

CANDLE Tutorial: Library Overview

ECP Annual Meeting, Houston, TX
January 2019

Jamaludin Mohd Yusof

Advanced Architectures and Applications
Los Alamos National Laboratory
jamal@lanl.gov



EXASCALE COMPUTING PROJECT



U.S. DEPARTMENT OF
ENERGY

Office of
Science

www.ExascaleProject.org

Talk Outline

- Introduction
- Benefits
- Library Overview
- Example Benchmark Workflow
- Simple Parameter Sweeps

Introduction

- Purpose
 - To streamline the writing of CANDLE-compliant codes
 - Allow rapid prototyping and exploration of hyperparameters
 - Integrate with the Supervisor framework
- Historical perspective
 - Consolidation of frequently used functionality from the Benchmark codes
 - Evolving to incorporate new functionality as needed
 - Improved usability over time

The CANDLE Environment

Hyperparameter Sweeps,
Data Management (e.g. DIGITS, Swift, etc.)

Workflow

Network description, Execution scripting API
(e.g. Keras, Mocha)

Scripting

Tensor/Graph Execution Engine
(e.g. Theano, TensorFlow, LBANN-LL, etc.)

Engine

Architecture Specific Optimization Layer
(e.g. cuDNN, MKL-DNN, etc.)

Optimization

Benefits provided by CANDLE

▪ Consistent

- Standardized network specification with a “default_model_file”
- Standardized command line intercept protocol
- Standardized default values across frameworks
- Ideal for testing the same problems with consistency on new DOE hardware

▪ Convenient

- Pass arguments via command line
 - Standard keywords parsed automatically, user can add new ones
- Modify the default file
 - Provide a new default model specification ‘--config_file new_default_model.txt’

Benefits provided by CANDLE

- Provides various utility packages that promote reuse and streamline code development
 - Actively developed, new functionality based on need
- Provides the pathway for inferencing, data-parallelism, automated sweeps of hyperparameters
- Availability of a robust framework for documentation and testing
- Pre-existing for containers such as Singularity (Ex. machines such as Theta, Titan, Cori, summitdev)
- Documentation: <https://ecp-candle.github.io/Candle/html/index.html>

Library Overview

- Integrated into the scripting level of CANDLE stack
 - Keras-based (for now)
 - Allows multiple backends (Tensorflow, CNTK, etc)
- Provides a single namespace for inclusion of useful functions into Benchmark codes
- Allows developers to decide which functions are exposed to users
 - `candle_keras` directory with `__init__.py` file sets included functions
- Allows reuse of non-Keras specific functions to create libraries for other languages

Library organization

- Utilities are organized by functionality
 - Transparent to the user, mostly framework agnostic
 - default_utils: create, modify parameter dictionary
 - file_utils: fetch and unpack data files
 - data_utils: load and manipulate data, enable UQ (via uq_utils)
 - generic_utils: callback function, standardize screen output
 - keras_utils: translation from CANDLER keywords, enhance Keras functionality
 - solr_keras: database functionality
 - viz_utils: visualize networks and output (prototype)

Example Benchmark Workflow

- `initialize_parameters`:
 - Read default values for the model from `default_model.txt` file
 - Automatically parse many standard ML hyperparameters
 - `learning_rate`, `batch_size`, etc
 - `data_url`, `train_data`, `test_data`...
 - Allows user to add other keywords via `additional_definitions`
 - Sets undefined hyperparameters to the corresponding Keras defaults to ensure consistency across backend frameworks
 - Early work revealed some performance differences were due to mismatched default settings

Example Benchmark Workflow

- Data management

- fetch_data

- Check if local data copy exists; if not, fetch data from data_url
 - allows separate train_data and test_data files, check MD5 hash if provided, untar if needed

- load_data

- Load labeled (xy) or unlabeled (x) data, perform various manipulations
 - Shuffle, scale, split into {train, validation, test}
 - Takes UQ operations to provide repeatable cross-validation data splits

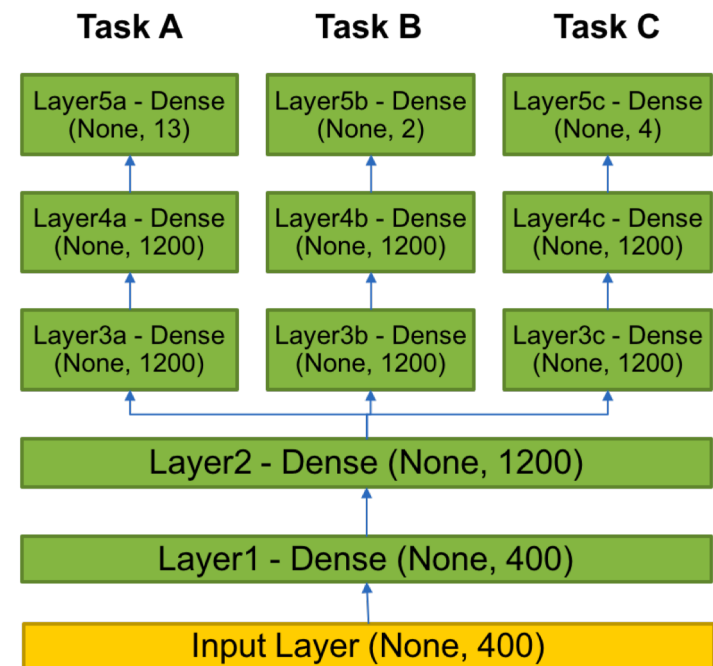
Example Benchmark Workflow

- Model build

- `keras_utils` translates CANDLE keywords into Keras methods
 - e.g. optimizer, initializer specification
- Extends Keras functionality where useful
 - `PermanentDropout`
- `default_utils` adds default `EXP000` and `RUN000` directory structure
- `solr_keras` adds monitoring and logging
 - Timeout functionality to respect job limits

Example benchmark

- **P3B1: Multi-task Deep Neural Net (DNN) for data extraction from clinical reports**
- **Overview:** Given a corpus of patient-level clinical reports, build a deep learning network that can simultaneously identify: (i) b tumor sites, (ii) t tumor laterality, and (iii) g clinical grade of tumors.
- **Relationship to core problem:** Instead of training individual deep learning networks for individual machine learning tasks, Build a multi-task DNN that can exploit task-relatedness to simultaneously learn multiple concepts.
- **Expected outcome:** Multi-task DNN that trains on same corpus and can automatically classify across three related tasks.



Original code

```
shared_nnet_spec= [ 1200 ]
individual_nnet_spec0= [ 1200, 1200 ]
individual_nnet_spec1= [ 1200, 1200 ]
individual_nnet_spec2= [ 1200, 1200 ]
individual_nnet_spec = [ individual_nnet_spec0, individual_nnet_spec1,
individual_nnet_spec2 ]

learning_rate = 0.01
batch_size = 10
n_epochs = 10
dropout = 0.0

## Read files
from data_utils import get_file
origin = 'http://ftp.mcs.anl.gov/pub/candle/public/benchmarks/P3B1/P3B1_data.tgz'
data_loc = get_file('P3B1_data.tgz', origin, untar=True, md5_hash=None,
cache_subdir='P3B1')
```

CANDLE code

```
gParameters = initialize_parameters()
...

# input layer
layer = Input( shape = ( input_dim, ), name= 'input' )
shared_layers.append( layer )

# shared layers
for k in range( len( shared_nnet_spec ) ):
    layer = Dense( shared_nnet_spec[ k ], activation=gParameters['activation'],
                  name= 'shared_layer_' + str( k ) )( shared_layers[ -1 ] )

    if gParameters['drop'] > 0:
        layer = Dropout( gParameters['drop'] )( shared_layers[ -1 ] )
    shared_layers.append( layer )

# individual layers
.....
```

CANDLE model file

```
data_url = 'ftp://ftp.mcs.anl.gov/pub/candle/public/benchmarks/P3B1/'
train_data = 'P3B1_data.tar.gz'
model_name = 'p3b1'
learning_rate = 0.01
batch_size = 10
epochs = 10
drop = 0.0
activation = 'relu'
out_activation = 'softmax'
loss = 'categorical_crossentropy'
optimizer = 'sgd'
metrics = 'accuracy'
n_fold = 1
shared_nnet_spec = '1200'
ind_nnet_spec = '1200, 1200:1200, 1200:1200, 1200'
feature_names = 'Primary site:Tumor laterality:Histological grade'
timeout = 1800
scaling = 'none'
output_dir = '.'
initialization='glorot_uniform'
```

Simple parameter exploration

- Provide a new default model specification
 - ‘--config_file new_default_model.txt’
- Overwrite individual parameters in the default model
 - ‘--learning_rate 0.1 -drop 0.1’
- Provides an easy way to perform individual experiments to probe the hyperparameter space

```
python myDNN.py -learning_rate 0.01 -run_id "run1"
```

```
python myDNN.py -learning_rate 0.02 -run_id "run2"
```
- Provides the pathway for automated sweeps of hyperparameters
 - → Supervisor workflows

Acknowledgments

This research was supported by the Exascale Computing Project (ECP), Project Number: 17-SC-20-SC, a collaborative effort of two DOE organizations - the Office of Science and the National Nuclear Security Administration, responsible for the planning and preparation of a capable exascale ecosystem, including software, applications, hardware, advanced system engineering and early testbed platforms, to support the nation's exascale computing imperative.



EXASCALE COMPUTING PROJECT



U.S. DEPARTMENT OF
ENERGY

Office of
Science

www.ExascaleProject.org