# Espam-AI manual

**Requirements & Installation:**
See README.md at the espamAI root folder or at https://git.liacs.nl/lerc/espam

**Execution from command line (after the tool is configured):**
1. cd <project root directory>
2. cd bin
3. ./espam [options]

**Execution from executable jar (after espam.jar is made by make jar-exec command):**
1. cd <jar directory>
2. java -jar espam.jar [options]

**Input:**
- Path to single DNN model in .onnx or .json format + options (see below) *or*
- Path to single CSDF model in .json format *or*
- Directory with several DNN models (not recomended for large DNN models)

**Output:**
- DNN model evaluation results in .json format printed to console *or/and*
- Files, generated from DNN model and intended for CSDF-models for
  CSDF-models processing tools such as Sesame/DARTS/SDF3

**Possibilities**:
- evaluation : Evaluation of DNN/CSDF model in terms of power/performance
- generation : Generation of files intended for CSDF-models
  processing tools such as Sesame/DARTS/SDF3\

**Tool running progress:**
(1) *Model(s) reading*: read input model(s)

(2) *Model(s) conversion*:
- conversion of input .onnx/json DNN model(s) to internal Network model(s)
  (initiated automatically after onnx/json model reading)
- conversion of internal Network model(s) to CSDF model(s):
  (initiated automatically, if evaluation flag if set or/and
   one or multiple *-csdf output files generation flags are set)

(3) *Output files generation [optional]:* generates for input model output files,
    according to set files generation options.

(4) *Model evaluation [optional]:* evaluates one or several input models
    in terms of power/performance, by means of the DARTS/SDF3 tool.
    TODO direct interface to SDF3 tool is not presented for current moment

## Command-line options

Command-line options that take arguments

| option | abbr | arguments | example |
|---|---|---|---|
| **General options** | | | |
| --generate | -g | Path to DNN/CSDFG model | --generate ./tests/lenet.json |
| --evaluate | -e | Path to DNN/CSDFG model | --evaluate ./tests/lenet.json |
| --consistency | -c | Path to DNN/CSDFG model | --consistency ./tests/lenet.json |
| **Model representation options (by default --layer-based option is set)** | | | |
| --layer-based | -lb | none | -lb |
| --neuron-based | -nb | none | -nb |
| --block-based | -bb | number of blocks | --bb 20 |
| --split-step [optional], for -bb models only | none | Number of child nodes, obtainbed after one layer splitting | --split-step 4 |
| **Hardware specification options** | | | |
| --time-spec | none | Path to input specification of operators times | --time-spec ./time_spec.json |
| --energy-spec | | Path to input specification of model energy | --energy-spec.json ./energy_spec.json |

Command-line flags (The command-line options that are either present or not).

| flag | abbr | description |
|---|---|---|
| **General flags** | | |
| --help | -h | Printout help |
| --version | -v | Printout program version |
| --verbose | -V | Printout program progress information |
| --copyright | none | Printout copyright |
| Input model type flag (--in-dnn is set by default) | | |
| --in-dnn | none | Input model is dnn model in .onnx or .json format |
| --in-csdf | none | Input model is csdf model in .json format |
| **Files generation flags** | | |
| --json | none | Generate DNN graph in .json format |
| --dot | none | Generate DNN graph image in .dot format |
| --json-csdf | none | Generate CSDF graph in .json format for DARTS |
| --xml-csdf | none | Generate CSDF graph in .xml fomat for SDF3 |
| --dot-csdf | none | Generate CSDF graph image in .dot format |
| --sesame | none | Generate code templates for Sesame |

| --wcet | none | Generate worst-case times specification for an input model |
|---|---|---|
| --wcenergy | none | Generate worst-case energy specification for an input model |
| --multiple-models | -m | Process not a single model, but multiple models. In this case path after –generate or –evaluate general flag should be the directory with several models. *Note: not recomended, escpecially for large models* |

## Files generation examples
For more information about files generation see comments, output files.

## Example 1

$ ./espam --generate ../src/espam/examples/DNN/json/lenet_LB.json –sesame -lb

*Description:*
      Generate layer-based CSDF graph for ../src/espam/examples/DNN/json/lenet_LB.json DNN model. Provide it with Sesame application.

*Expected output:* files folder
&lt;lenet_LB&gt;
      -app //sesame application in layer-based mode

## Example 2
$ ./espam --generate ../src/espam/examples/CSDF/json/tinyCSDF.json --in-csdf --dot-csdf --sesame

*Description:*
      For ../src/espam/examples/CSDF/json/tinyCSDF.json CSDF graph generate .dot file and Sesame application.

*Expected output:* files folder
&lt;tinyCSDF&gt;
      -app //sesame application for an CSDF graph
      -sdfg //CSDFG files directory
          -dot //contains tinyCSDF.dot file of CSDF model

## Example 3
$ ./espam --generate ../src/espam/examples/DNN/json/lenet_NB.json --sesame --dot --json-csdf --xml-csdf --dot-csdf -nb

*Description:*
      Generate neuron-based CSDF graph for ./tests/json/lenet_NB.json DNN model. Provide it with Sesame application, .dot files of DNN and CSDF model, .json and .xml files of CSDF model.

*Expected output:* files folder
&lt;lenet_NB&gt;
      -app //contains sesame application in -nb mode
      -dot //contains  lenet_NB.dot file of DNN model
      -sdfg //corresponding CSDF graph directory
          -dot //contains lenet_NB.dot file CSDF model
          -json //contains lenet_NB.json CSDF model suitable for DARTS\
          -xml //contains lenet_NB.xml CSDF model suitable for SDF3

**Example 4**

$ ./espam --generate ../src/espam/examples/DNN/onnx/mnist.onnx --sesame --dot --json-csdf --xml-csdf --dot-csdf -bb -100 --split-step 4

*Description:*

       Generate block-based CSDF graph for ./tests/onnx/mnist.onnx DNN model. Split model until it reaches <=100 blocks with –split-step = 4 (For more information about layers splitting see Cooments, block-based models). Provide it with Sesame application, .dot files of DNN and CSDF model, .json and .xml files of CSDF model.

*Expected output:* files folder
<CNTKGraph_NB>
       -app //contains sesame application in -bb mode
       -dot //contains  CNTKGraph.dot file of DNN model
       -sdfg //corresponding CSDF graph directory
          -dot //contains CNTKGraph.dot file of CSDF model
          -json //contains CNTKGraph.json CSDF model suitable for DARTS\
          -xml //contains CNTKGraph xml CSDF model suitable for SDF3

## Model evaluation examples

### Example 1

$  ./espam --evaluate ../src/espam/examples/DNN/json/lenet_LB.json --lb

*Description:*

       Evaluate ../src/espam/examples/DNN/json/lenet_LB.json DNN model in terms of power/performance in -lb mode with default time specification.
*\* For more details about time and energy specification see 'Comments. hardware specification options' chapter below.*

*Expected output:*
    { "id": 0, "execution_time": 4361469.0, "energy": 5.8683824566753834E-9, "memory": 278280.0, "processors": 2 }

### Example 2

$  ./espam --evaluate ../src/espam/examples/CSDF/json/tinyCSDF.json --in-csdf

*Description:*

       Evaluate ../src/espam/examples/CSDF/json/tinyCSDF.json CSDF model in terms of power/performance with default time and energy specification.

*Expected output:*
    { "id": 0, "execution_time": 18.0, "energy": 6.4387499999999995E-9, "memory": 96.0, "processors": 3 }

### Example 3

$  ./espam --evaluate ../src/espam/examples/DNN/json/lenet_NB.json -nb --time-spec
    ../src/espam/examples/time_spec/wcet.json --energy-spec
    ../src/espam/examples/energy_spec/wcenergy.json

*Description:*

Evaluate ../src/espam/examples/DNN/json/lenet_NB.json DNN model in terms of power/performance in -nb mode with time specification given in wcet.json file and energy specification given in wcenergy.json file.

*Expected output:*
> { "id": 0, "execution_time": 359499.0, "energy": 42.371048349077284, "memory": 787744.0, "processors": 14 }

## Example 4

$ ./espam --evaluate ../src/espam/examples/DNN/onnx/mnist.onnx --sesame -bb -100 --split-step 4

*Description:*
> Evaluate ./src/espam/examples/DNN/onnx/mnist.onnx DNN model in terms of power/performance. Split model until it reaches <=100 blocks with –split-step = 4 (For more information about layers splitting see Cooments, block-based models) with default time and energy specifications.

*Expected output:*
> { "id": 0, "execution_time": 9932160.0, "energy": 4.460457589285961E-9, "memory": 83424.0, "processors": 2 }

## Consistency checkout examples

## Example 1

$ ./espam --consistency ../src/espam/examples/DNN/onnx/mnist.onnx

*Description:* Check if ../src/espam/examples/DNN/onnx/mnist.onnx model is consistent.

*Expected output:*
> input dnn consistency: true

## Example 2

$ ./espam --consistency ../src/espam/examples/DNN/onnx/mnist.onnx

*Description:* Check if ../src/espam/examples/DNN/onnx/mnist.onnx model is consistent.

*Expected output:*
> input dnn consistency: true

## Comments

**1. Output files**
> Output files are generated in output models directory, specified during espamAI configuration. By default output files are generated in espamAI/execution/directory/output_models. For every model created a directory, named after the model. E.g. for DNN called "LeNet", the LeNet directory will be created.

Inside the folder there are files generated in accordance with the file generation flags. For all possible files the folder will have the following structure (each file or direcory is optional)

The files structure:
<Model_name>
> - app // directory with sesame application
> - Model_name_wcet_spec.json //time specification in .json format
> - energy_spec.json //time specification in .json format
> - dot //directory with DNN model in .dot format
> - json //directory with DNN model in .json format
> - sdfg //corresponding CSDF graph directory

- dot //directory with CSDF model in .dot format
- json //directory with CSDF model in .json format for DARTS
- xml //directory with CSDF model in .xml format for SDF3

*NOTE: if directory with this name already exist it will be overwritten!*

## 2. Neuron/Layer and Block-based DNN models

EspamAI provides several types of DNN--> CSDF model granularity.
For layer-based model each node of a CSDF graph represents one layer of an input DNN model.
For neuron-based model each node of a CSDF graph represents one neuron of an input DNN model.
For block-based model each node of a CSDF graph represents one block of an input DNN model,
where block is an abstraction of DNN model layer and DNN model neuron.
Block-based model takes on input layer-based DNN model and transforms it according to the
following algorithm:

**While** (number of blocks<expected  && split_up flag = true):
    Layer bottleneck = find_bottleneck(); // bottleneck node search is preformed by DARTS
    **if** (bottleneck can be split up)
        evaluate number of layers after bottleneck layer splitting;
        **if** (number of layers after splitting>expected)
            **return**;
        Split bottleneck into –split-step child Layers;
        Split all Dependent layers (nonlinear/maxpool ones), following the bottleneck layer
    **else**
        **return**;

As transformation is performed over DNN models,files (.json / .dot etc) generated from the model
in -bb mode will be different from ones for input (layer-based) DNN model. If .json for block-based
model is already generated, it can be reused without additional splitting in -lb mode.

## 3. Hardware specification options

EspamAI provides a possibility to manage time and energy specifications, used for input
models evaluation.

### 3.1. Time specification

Time specification is provided as a list of k-v pairs *operation: execution time* and stored in .json
format. The operation ia s single opetator, performed in an CSDF node such as Convolution or
Subsampling.
There is a list of parametrized operations, supported by default. The default supported operators list
is provided in ../espam/src/espam/docs/espamAI/DNN_supported_operators.pdf

Example of specification with only default operators is provided below:
```
{
  "READ": 1,
  "WRITE": 1,
  "MAXPOOL": 2,
  "CONV": 3,
  "SOFTMAX": 1,
  "NONE": 0
}
```
Each default operator have a number of parameters and 'time per pixel' value. The parametrized
operator time is a funcion of 'time per pixel' and actual parameters, automatically extracted from
corresponding CSDF node during time evaluation.

The specification can be extended by non-parametrized or arbitrary operators:

```
{
  "READ": 1,
  "WRITE": 1,
  "MAXPOOL": 2,
  "CONV(5_5)": 75,
  "CONV": 3,
  "DENSEBLOCK_NONE(1_400,400_120)": 10400,
  "WRITE": 1,
  "SOFTMAX": 1
}
```

The non-parametrized operator name will have a priority for time evaluation, e.g. for the specification above for Convolution 5x5 operator, the "CONV(5_5)": 75 k-v pair will be selected for 5x5 Convolution. For Conv mxn, where m!=5 and n!=5 the "CONV": 3 k-v pair will be used.

It might happen that an arbitrary operator, which is not mentioned in supported list is used. In this case the operator time will be set to 1 and the warning *<operator> unknown execution time. Default time = 1 is set for <operator>* will be given.

**Example**

To set up your own time specification, perform 3 following steps:

1. Generate current specification for an input model:

        $ ./espam --generate ../src/espam/examples/DNN/json/lenet_NB.json -nb –wcet

*Expected output:*

        File with default times: ../output_models/lenet_LB/lenet_NB_wcet_spec.json of structure:

```
{
  "INPUT": 1,
  "OUTPUT": 1,
  "MAXPOOL": 2,
  "CONV(5_5)": 75,
  "CONV": 3,
  "DENSEBLOCK_NONE(1_400,400_120)": 10400,
  "WRITE": 1,
  "SOFTMAX": 1,
  "READ": 1,
  "MAXPOOL(2_2_10)": 80,
  "DENSEBLOCK_NONE(1_120,120_84)": 2160,
  "DENSEBLOCK_SOFTMAX(1_84,84_10)": 504,
  "ReLU": 1,
  "CONV(5_5_6)": 450,
  "MAXPOOL(2_2_28)": 224,
  "NONE": 0
}
```

2. Change the times .json file manually (e.g. set "CONV(5_5)": 95) in any text editor
3. Call the model with changed specification:

        $ ./espam --evaluate ../src/espam/examples/DNN/json/lenet_NB.json -nb --time-spec
../output_models/lenet_NB/lenet_NB_wcet_spec.json

*Expected output:*

{ "id": 0, "execution_time": 359499.0, "energy": 4.237104834907729E-8, "memory": 787744.0, "processors": 14 }

### 3.1. Energy specification

Energy model is based on based on:

Di Liu, Jelena Spasic, Gang Chen, and Todor Stefanov,
Energy-Efficient Mapping of Real-Time Streaming Applications on Cluster Heterogeneous MPSoCs", In Proc. "13th Int. IEEE Symposium on Embedded Systems for Real-Time Multimedia (ESTIMedia'15), pp. 1-10, Amsterdam, The Netherlands, Oct. 8-9, 2015.

### Example

To set up your own energyspecification, perform 3 following steps:
1. Generate current specification for an input model:

  $ ./espam --generate ../src/espam/examples/DNN/json/lenet_NB.json -nb –wcenergy

*Expected output:*

  File with default worst-case energy: ../output_models/lenet_LB/energy_spec.json of structure:

```
{
  "alpha": 3.03E-9,
  "beta": 0.155,
  "b": 2.621
}
```

2. Change the energy specification in any text editor (e.g. set beta to 0.275)

*IMPORTANT: Energy specification should always contain alpha, beta and b parameters, only values can be changed!*

3. Call the model with changed specification:

  $ ./espam --evaluate ../src/espam/examples/DNN/json/lenet_NB.json -nb  --energy-spec
    ../output_models/lenet_NB/energy_spec.json

*Expected output:*

  { "id": 0, "execution_time": 359499.0, "energy": 4.237104834907729E-8, "memory": 787744.0, "processors": 14 }

## 4. Simulate model running with Sesame

To run model, generated with –sesame espamAI generation flag (see commands), you will need Sesame simulation tool installed.

With installed Sesame simulation tool:
1. Go to sesame root
2. Set up sesame environment
3. Go to directory with an application, generated by espamAI
4. Provide sesame mapping file or reate virtual mapping for your application
5. Go to application sources.
6. Generate .o files for every CSDF node and .so library, contains all the application + simulate application running.
7. Go back to application root folder
8.Create text files with application processes traces.

### Example

An example for todorsNet_NB Sesame application with virtual mapping:
1. $ cd ../sesame/

2. $ source sesame.env

3. $ cd .../todorsNet_NB/

4. $ AppVirtualMapGenerator app/todorsNet_NB_app.yml > todorsNet_NB_appvirt_map.yml

5. $ cd app

6. $ make runtrace

7. $ cd ../

8. $ for i in `ls trace*`; do traceprinter $i > $i.txt; done

*Expected output:*

      Run traces for every CSDF node with sequences of read/execute/write primitives, corresponding to the CSDF graph functionality.


## 5. Generate images from .dot files

To generate images from .dot files use graphviz https://www.graphviz.org/

It should installed on Linux by default. In case it is not installed, use

                      **$ sudo apt-get install graphviz**

                                **Example**

To generate .png image of todorsNet_NB.dot file use

      $ dot -Tpng ../espam/output_models/todorsNet_NB/dot/todorsNet_NB.dot -o
      ../espam/output_models/todorsNet_NB/dot/todorsNet_NB.png

For more output file formats see

      $ dot –help


## 6. Recomendations and possible issues.

3.1. It is recomended to run tests after the tool is build.

TODO provide run tests after tool installation

3.2. Tool is not accounted on huge graph models, so it is not recomended to use -nb mode for large DNNs. Otherwise, tool may work really long or even fall with memory (JavaHeapSpace/Python) errors.

For the same reason, it is not recommended to use –generation option with too many flags and multiple-models processing (--evaluation or ---generation) for heavy models, especially in -nb mode

3.3. Do not forget to set –in-csdf flag for input CSDF graph models, otherwise CSDF model will not be processed.

3.4. If you are not sure about input model consistency, use –consistency option to check it.

3.5. Tool may complain on 'unsafe' libraries such as protobuf-java.jar.

3.6. If tool is frozen, kill the corresponding command line process and try to use the same command with –verbose (v) option. It will be shown on which step (model reading, model conversion, model evaluation etc.) the problem occured and more details will be given.

3.7. If CSDF model is sent on input (with flag –in-csdf) , no specific DNN features (such as weights and proper data formats) will be taken into account. Thus, for processing DNN models it is recomended to use DNN .json/.onnx files. JSON files might be even preferable due to their small size.


## 7. References

1. EspamAI https://git.liacs.nl/lerc/espam/tree/espamAI
2. DARTS tool http://daedalus.liacs.nl/daedalus-rt.html

3. Energy model Di Liu, Jelena Spasic, Gang Chen, and Todor Stefanov, Energy-Efficient Mapping of Real-Time Streaming Applications on Cluster Heterogeneous MPSoCs", In Proc. "13th Int. IEEE Symposium on Embedded Systems for Real-Time Multimedia (ESTIMedia'15), pp. 1-10, Amsterdam, The Netherlands, Oct. 8-9, 2015.
Electronic version is available on http://liacs.leidenuniv.nl/~stefanovtp/publications.html
4. Sesame tool: http://sesamesim.sourceforge.net/