

# *expLanes: beautiful computational experiments*

Mathieu Lagrange

December 9, 2015

<i>The scientific method</i>	1	<i>Architecture of an experiment</i>	8
<i>About you</i>	2	<i>Configuration of an experiment</i>	10
<i>Features</i>	3	<i>Exposition of observations</i>	14
<i>Benefits</i>	4	<i>Publishing your experiment</i>	18
<i>Installation</i>	4	<i>Advanced usage</i>	20
<i>In a nutshell</i>	5	<i>Recommended readings</i>	22

## **The scientific method**

Since it was expressed in its modern form by Roger Bacon in the 13th century, the scientific method has demonstrated time and again its merit, as the process that underpinned all the discoveries and innovations upon which our modern world is built.

Sadly, modern ways of doing research impose strong pressure on the time and efforts that can be allocated to a project, and important steps of the scientific method are often neglected. We believe that this quest for speed adversely reduces the meaningfulness of the research results that are nowadays published.

Admittedly, following strictly the scientific method can be tedious and shortcuts might be tempting. *expLanes* is designed to assist you in the most tedious and most error prone steps, and will hopefully help you dedicate more time to the fundamental steps of the scientific method.

Quickly put, the scientific method can be divided into several steps, each of which may have to be iterated. On the following table, we stated where the exPlanes framework can be helpful during this iteration process.

The scientific method <sup>1</sup>		
Phases	Steps	exPlanes
<b>Analysis</b>	Describe problem	
	Set performance criteria	
	Investigate related work	
	State objective	
<b>Hypothesis</b>	Specify solution	
	Set goals	
	Define factors	+
<b>Synthesis</b>	Postulate performance metrics	+
	Implement solution	++
	Design experiments	+++
	Conduct experiments	++++
	Reduce results	+++
<b>Validation</b>	Compute performance	++
	Draw conclusions	+
	Prepare documentation	+
	Solicit peer review	

<sup>1</sup> P. Bock and B. Scheibe. *Getting it right: R & D methods for science and engineering*. Academic Press, 2001

## About you

Let's assume that you are practicing computational experiments. Most probably, this takes you time and a lot of care, and may be an endless source of frustration depending on whether you are:

- **a Master's student?** Then, you may at some point consider the fact that the problem is not simply to come with a new idea and implement it. To contribute significantly to the research community that you are striving to be part of, you need to compare your method with existing ones. This process is tedious, hard (if not impossible) and involves a lot of coding and knowledge about large scale data processing, statistical analysis and reporting of quantitative data.
- **a PhD student?** You have several years to dedicate to a research project. Doing it well will help you stay motivated and efficient. But how? Several years of work means a lot of code, a lot of bugs, a lot of failures and hopefully some gain of knowledge for you and your research community. How will you keep track of those many experiments? How will you efficiently document them? How will you quickly report your progress to the members of your research team? How will you publish your research in a reproducible way?
- **a Post doctoral fellow?** You are now an established researcher, with many ideas about what could be done in order increase

knowledge in your community. But you also have to juggle with the many different projects you are involved in. Keeping track of all those projects and being able to easily switch between them is mandatory for success. For example, being able to re-run years old experiments in order to efficiently satisfy a reviewer request is critical for your career.

- **a Full time researcher?** Besides research, you have many duties that shred the time you can allocate to pursuing your many personal research projects, and the time needed to switch between projects can easily become excessive with respect to your overall available time. More importantly, the free time you have is usually not in front of your desktop. Moreover, you are supervising several students whose research work, as valuable as it may be, is more often than not plagued with such a level of organizational idiosyncrasy as to make it close to useless for the next student working on the project—or yourself.

And for all or others, maybe you are heading towards sharing your code, but not quite confident with your programming expertise yet, and you do not have time to improve your experimental code into a sharable state<sup>2</sup>? Or maybe you need solid experimental data to pitch a new technology?

<sup>2</sup> <http://sciencecodemanifesto.org>

expLanes is specifically designed with those matters in mind, and will hopefully help you reduce the specific burden that keeps you from reaching this goal which is one of the most important steps towards the expansion of knowledge in science and engineering: reproducibility<sup>3</sup>.

<sup>3</sup> R. D. Peng. Reproducible research in computational science. *Science*, 334(6060):1226–1227, 2011; and P. Vandedewalle, J. Kovacevic, and M. Vetterli. Reproducible research in signal processing. *Signal Processing Magazine, IEEE*, 26(3):37–47, May 2009

## Features

expLanes is a software framework currently implemented in Matlab that

- provides a high level abstraction for running computational experiments,
- allows multiple users per experiment,
- allows multiple processing platforms to be used,
- allows a strong decoupling of the 3 major experimental phases:
  1. coding,
  2. processing,
  3. reduction of results.

Concretely, expLanes lets you identify the operational parameters (factors of the experiment) of a process or algorithm, and their ranges of variation (modalities of the experiment), then proceeds to automatically explore all or a subset of the space of possible parameter combinations, using a dataset and evaluation function of your

choice. Finally, it generates a report compiling the obtained results. This can for example be used to evaluate an average performance, tune an algorithm for a specific type of data, or confirm empirically that a process behaves satisfactorily in all situations.

## Benefits

1. The user can focus on solution code,
2. evaluated with standard experimental designs.
3. Bugs and the time needed to solve them is reduced,
4. context switching between projects is much easier,
5. as well as diffusion of reproducible code.

## Installation

explanes operates in the Matlab environment. A recent version of this software (> R2014) is needed as well as the Statistical toolbox. The Parallel toolbox is needed for in-session multi-core processing.

As much as possible, explanes uses built-in tools to manipulate data. However some functions are handled through the use of basic Unix tools:

1. **pdflatex**: handles compilation of reports written in  $\text{\LaTeX}$ .
2. **rsync**: allows backup, generation of bundles and data communication between hosts.
3. **ssh**: allows secure connection between hosts.
4. **screen**: allows handling of several detached terminals for background processing.

Please note that those tools must be available on every host where an explanes project is processed. That said, explanes does not enforce their availability. For example, without `pdflatex`, no reports will be generated.

Probing access to those tools by your explanes project (in this documentation, the project we will demonstrate is called `geometricShape`) can be done using the following command:

```
geometricShape('probe', 1);
```

It also checks if the different data paths and processing hosts are reachable. Each of those configuration parameters can be set for each user by editing the configuration file of the project. A default configuration file to be used as a template when creating a new project is available at `~/explanes/<userName>Config.txt`.

## Linux

On Linux hosts, the installation should be straightforward. For example, on Debian, Ubuntu systems, the following packages should be sufficient:

```
texlive-full rsync openssh-server openssh-client ...
screen
```

Installation may be achieved by typing in a terminal:

```
sudo apt-get install <packageNames>
```

Minimal configuration might be needed to configure ssh on client and server sides depending on your network configuration. In `expLanes` ease of use is achieved by using empty passphrase logging.

## Mac OS

Several alternatives are available to install Unix packages: Fink, MacPorts and Homebrew. You should avoid the use of `/usr` for installation, and the `PATH` environment variable should be properly set.

## Windows

For `pdflatex`, please use MiKTeX, for the remaining, please use the Cygwin environment. Some tweaks might be needed in order to ensure that the paths and usernames are correct.

## In a nutshell

`expLanes` is designed to provide you with as streamlined set of tools to efficiently build the computational environment you need in order to gain knowledge about a research statement.

In order to present the typical workflow of a `expLanes`-driven experiment, we detail in this section the steps to explore a trivial research statement: gaining knowledge about the base area and the volume of a few 3 dimensional geometrical shapes. For reference, this project is available in the `demonstrations` directory. Complete documentation of the commands is given in the next section.

First, we assume that the `expLanes` framework is in your path, if not, please type the following in your command window:

```
addpath(genpath('<pathToExplanes>'));;
```

### Create project

To create the `geometricShape` project, call:

```
expCreate('geometricShape');
```

The command ends by moving into the experiment directory.

This documentation can be loaded at any time using the help command:

```
geometricShape('h');
```

*Define steps*

The processing steps can be now be instantiated:

```
geometricShape('addStep', 'base');
geometricShape('addStep', 'space');
```

Alternatively, the 3 previous commands can be run at once:

```
expCreate('geometricShape', {'base', 'space'});
```

*Define factors*

We are interested in the potential impact of the different attributes (shape, color, radius, width, height) of the geometric shape on its base area and volume. For each of those attributes, we define a set of values that shall be explored during the experiment.

The shape can be a cylinder, a pyramid or a cube:

```
geometricShape('addFactor', ...
    {'shape', {'cylinder', 'pyramid', 'cube'}});
```

The shape can be blue or red:

```
geometricShape('addFactor', {'color', {'blue', 'red'}});
```

The radius of the cylinder (modality 1 of factor 1) can be 2, 4, or 6 meters:

```
geometricShape('addFactor', {'radius', '[2, 4, 6]', ...
    '', '1/1'});
```

The width of the pyramid and the cube (modalities 2 and 3 of factor 1) ranges from 1 to 3 meters:

```
geometricShape('addFactor', {'width', '1:3', '', ...
    '1/[2 3]});
```

The height of the shape is only relevant for processing step 2 and for the cylinder or the pyramid (modalities 2 and 3 of factor 1):

```
geometricShape('addFactor', {'height', '2:2:6', '2', ...
    '1/[1 2]});
```

Factors can be viewed by typing: `<experimentName>()`; , thus

`geometricShape()`; now returns:

Factors :

```
1   shape = = = {'cylinder', 'pyramid', 'cube'}
2   color = = = {'blue', 'red'}
3   radius = = 1/1 = [2, 4, 6]
4   width = = 1/[2 3] = 1:3
5   height = 2 = 1/[1 2] = 2:2:6
```

The factors can be edited in the file

`<shortExperimentName>Factors.txt`, that is `geshFactors.txt`.

A setting is a set of modalities, one of each factor of interest.

*Implement processing steps*

The first step is dedicated to the computation of the base area of the shape. The solution code is implemented in `gesh1base.m`, where we simulate that  $\pi$  is estimated with a given precision, but 100 measurements have been made:

```
uncertainty = randn(1, 100);
switch setting.shape
```

```

    case 'cylinder'
        baseArea = (pi+uncertainty)*setting.radius^2;
    otherwise
        baseArea = setting.width^2;
end
store.baseArea = baseArea;

```

The second step builds upon the result of the first processing step to compute the volume (implemented in `gesh2volume.m`):

```

switch setting.shape
    case 'cube'
        volume = data.baseArea*setting.width;
    case 'cylinder'
        volume = data.baseArea*setting.height;
    case 'pyramid'
        volume = data.baseArea*setting.height/3;
end

```

### Define observations

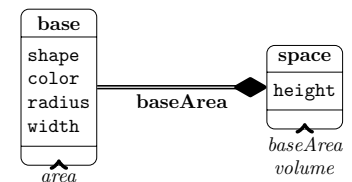
Observations for the 2 processing steps are the following. First step:

```
obs.area = baseArea;
```

Second step:

```
obs.baseArea = data.baseArea;
obs.volume = volume;
```

The command `geometricShape('f')` generates a diagram view of the experiment.



### Process

```
geometricShape('do', 1);
```

run the base step over all 18 settings.

```
geometricShape('do', 0, 'mask', {[1 2] 0 1});
```

runs successively each step over the cylinders of radius 2 and all the pyramids.

### Expose observations

Upon completion of the processing, the results of the last processing are displayed in the command window:

This display can be obtained by issuing the following command:

```
geometricShape('display', 2, 'expose', '>', ...
    'mask', {[1 2] 0 1});
```

The exposition of relevant observations is important for efficient computational experimentation. `explanes` provides many tools for this purpose, type `help expExpose` for quick reference.

The command:

```
geometricShape('display', 2, 'mask', {1 0 1}, ...
    'expose', {'t', 'obs', 3, 'sort', 1});
```

displays the volume (observation 3) of the step volume (2) for each cylinder of radius 2 sorted according to the first on a table (t). Red color indicates best performance, and blue ones, performances which have not been found statistically different from the best one.

	color	height	volume
1	blue	6	78.35 (27.90)
2	red	6	73.84 (22.61)
3	blue	4	52.23 (18.60)
4	red	4	49.23 (15.07)
5	blue	2	26.12 (9.30)
6	red	2	24.61 (7.54)

For this example, even if the blue cylinder is bigger than the red one due to the uncertainty in the estimation of  $\pi$ , this difference is not significant, so the blue and red cylinders shall be considered of equivalent volume.

## Report

The file `<shortExperimentName>Report.m`, that is `geshReport.m` in this example, allows you to generate a report that compiles several expositions of observations:  $\text{\LaTeX}$ table (l), bar plot (b), and box plot (x).

```
config = expExpose(config, 'l', 'obs', 3, ...
    'mask', {1 0 1}, 'save', 'geol');
config = expExpose(config, 'b', 'obs', 3, ...
    'mask', {1 0 1}, 'save', 'geob', 'orientation', ...
    'h');
config = expExpose(config, 'x', 'obs', 3, ...
    'mask', {1 0 1}, 'save', 'geox', 'orientation', ...
    'h');
```

Several reports can be handled for the same experiment using the key `'reportName'`, and if the name of the report contain the key `'Slides'`, a slide presentation layout is used. For example, the command:

```
geometricShape('report', 'rcv', ...
    'reportName', 'presentationSlides');
```

generates a report with base name `presentationSlides` in a slides presentation layout. For debug information about the compilation, please add the `'b'` flag to the `report` value.

## Demonstrations

Together with the examples detailed in this manual, some demonstrations are available in the `demo` directory. At the root of the directory, some reports are given as examples of what can be generated within an `explanes` experiment.

Those reports are available online:

1. **Computation of distances:** <https://github.com/mathieulagrange/explanes/blob/master/demo/distanceComputation.pdf>
2. **Clustering:** <https://github.com/mathieulagrange/explanes/blob/master/demo/clusteringData.pdf>
3. **Classification:** <https://github.com/mathieulagrange/explanes/blob/master/demo/classifyData.pdf>
4. **Audio source separation:** <https://github.com/mathieulagrange/explanes/blob/master/demo/denoiseAudio.pdf>

## Architecture of an experiment

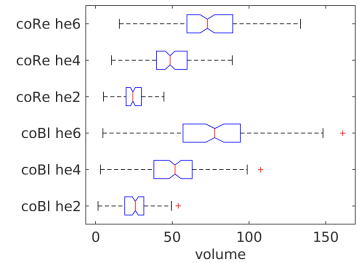
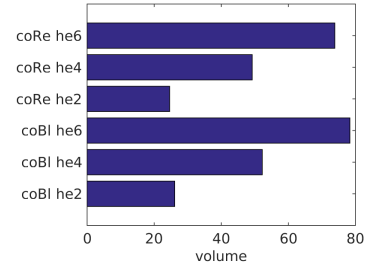
An `explanes` experiment has a specific directory architecture.

### Code

The main directory hosts the processing routines (named `codePath` in your configuration). It has the following files:

Table 1:  $\text{\LaTeX}$  table output.

color	height	volume
blue	2	$26.12 \pm 9.30$
blue	4	$52.23 \pm 18.60$
blue	6	<b><math>78.35 \pm 27.90</math></b>
red	2	$24.61 \pm 7.54$
red	4	$49.23 \pm 15.07$
red	6	<b><math>73.84 \pm 22.61</math></b>





- `<projecName>`: main entry point of the experiment
- `Factors`: text file encoding the factors of the experiment
- `Init`: processed before any processing step
- `Steps`: implement each processing units
- `Report`: handle the generation of report that compile several types of expositions of observations

### Configuration

The processing environment of the experiment can be controlled with the files hosted in the `config` directory.

Almost every processing unit of an `explanes` experiment has access to a structure named `config`, which is imported from a file that is specific to each user. This file is generated at the creation of the experiment from a user specific template file that can be found in `~/explanes`.

The syntax of this file allows the user to handle several configurations depending upon which machine the experiment is processed on. For example, let us assume that the project handles 4 machines with different file architectures:

```
machineNames = {'toto', 'yoyo', 'dodo', 'momo'}
codePath = {'~/code/geometricShape', ...
            '/lab/code/geometricShape'}
```

If the experiment is processed on the machine `'toto'`, or `'yoyo'`, the `codePath` in the `config` structure is equal to `'~/code/geometricShape'`.

If the `codePath` has fewer entries than the one dedicated to the names of the machines, the last entry is chosen. Thus, if the experiment is processed on the machine `'dodo'`, or `'momo'`, the second entry is chosen, that is `'/lab/code/geometricShape'`.

Every entry of the configuration file can be overwritten at the command line call. For example:

```
geometricShape('do', 0, 'sendMail', 1)
```

processes every steps and sends an email at the end of the processing.

New entries can also be added and accessed in the processing files of the experiment for specific user usage:

```
geometricShape('myEntry', 'toto')
```

outputs:

Warning. The command line parameter `myEntry` is not found in the `Config` file. Setting anyway.

### Data

Access to input data is handled in the configuration file with the entry `inputPath`. Processing data is handled with the entry `dataPath`. Before processing any step, a directory is created with the name of the step, and all processing data generated by this step is stored here.

Observations are usually of smaller size. As such it may be useful to store them in a different location such as a cloud filesystem. The `obsPath` stores the observation data in the same directory architecture. If left empty, the `obsPath` is set equal to `dataPath`.

### Report

The report directory stores expositions of observations in many flavors.

- **report/**: contains editable  $\text{\LaTeX}$  files, one for each `reportName` (defaults to `<projectName>.tex`) and a BibTeX file that can be edited and processed in a standard fashion. Expositions of observations are added to the  $\text{\LaTeX}$  file prior to compilation at the location of the flag `expLanesInsertionFlag`.
- **report/figures**: when the `'save'` parameter is added to a visual exposition, the resulting figure is stored in several formats (`.fig` Matlab Figure, `.eps`, `.png`, and `.pdf`). The numerical data is also made available in a `.mat` file.
- **report/reports**: contains the generated reports with the following naming template: `<reportName><userName><date>.pdf`. The data of each exposition of the corresponding report is made available in a `.mat` file.
- **report/tables**: when the `'save'` parameter is added to a table exposition (`'t'` or `'l'`), the resulting table is stored in several formats ( $\text{\LaTeX}$  floating table,  $\text{\LaTeX}$  tabulary array, `.csv` file). The numerical data is also made available in a `.mat` file.
- **report/tex**: this directory is used for internal `expLanes` usage. Please do not edit any files in it, as they may be deleted.

## Configuration of an experiment

An `expLanes` experiment can be controlled by editing the configuration file or overloading it with command line settings. Each of those parameters can further be accessed at run time through the variable `config`.

### User

- **emailAddress**: email address used to acknowledge when processing is started or over server side (leave empty for no mailing service), can be a cell array of email addresses
- **completeName**: complete name of user

### Project

- **probe**: probe tools (1), paths (2), hosts (3), or all (0)

- **generateRootFile**: regenerate the root file of the experiment (newer versions of the expLanes lib may require it)
- **addFactor**: add a factor to the experiment. Format of the request: {'name', 'modalities', 'steps', 'selector', defaultModalityId, rankInFactorFile}
- **removeFactor**: remove a factor from the experiment. Format of the request: {rankInFactorFile, defaultModalityId}
- **addStep**: add a step to the experiment. Format of the request: 'name' (last step rank) or {'name', rank}
- **removeStep**: remove a given step from the experiment. Format of the request: rank
- **readMe**: fill the README.txt file with information about replication of the experiment (boolean)

### *Hostnames and Paths*

- **machineNames**: names of the machines as a cell array of 1) strings or 2) cell arrays of strings
- **inputPath**: path to input data (usually accessed read only)
- **codePath**: path to the code repository of the experiment
- **dataPath**: path to the data repository of the experiment (structured with one directory per step)
- **obsPath**: path to the data repository of the experiment for observation data (same as dataPath if left blank)
- **backupPath**: path to backup of transferred data
- **exportPath**: path to export replication of the experiment
- **matlabPath**: path to the directory of the Matlab binary
- **toolPath**: path to Unix tools (pdflatex, rsync, ssh, screen)

Every path can be a cell array of string. In this case, the path is selected in the cell array according to the rank of the host in the machineNames field.

### *Code*

- **dependencies**: code dependencies to load and send server side as a cell array of strings with paths to the dependencies
- **localDependencies**: dependencies (including expLanes) can be part of the experiment in a dependencies directory: 0 do not use local versions, 1 use local versions, 2 update local versions
- **codeVersion**: version tag of the code of the experiment

### Computing

- **do**: processing steps to execute: -1 none, 0 all, >0 processing step by numeric id
- **resume**: do not perform computation if data and observation files with runID>resume are already there
- **parallel**: use parallel processing; if |parallel|>1 specify the number of cores, if array set at step level, if >0 parallelize settings, if <0 parallelize within each setting (the code of the step shall use parfor or the like), can be a numeric value for all processing steps or an array of numeric value, one per step
- **mask**: factor mask: cell array defining the modalities to be set for the factors (0 do all, 1 first modality, [1 3] first and third modality), can be a cell array of cell arrays
- **design**: use a canonical experimental design: assume cell array of {[factors as numericIds], number of values per factor (0 complete set of values), type of plan as string 'f' (factorial) or 'o' (one factor at a time), [seed as a vector of index of factor]}
- **dummy**: dummy mode: allow for short computations to dry run the experiment. Set as a numeric value. 0: no dummy mode, >1 generate stored data and observations flagged with the numeric value
- **setRandomSeed**: set the random seed at init for replicability purposes, 0 do not set, >0 set to value.
- **host**: host index to run the experiment: 0: seek by hostName, >0: server mode, <0: local mode, 2.1 means the first host of the second compound (cell array)
- **log**: log level (set to 0 for no log)
- **progress**: show processing progress 0: none, 1: graphical bar if on local mode, 2: verbose output, 3: terse output
- **exitMatlab**: exit matlab at the end of the computation
- **recordTiming**: show timing observations

### Data

- **encodingVersion**: encoding type of mat files (doc save for specs)
- **store**: perform computation to be stored: (1, 0) load data of previous step, (-1, 0) load data result of the current step
- **retrieve**: retrieve needed data server side: -1 no retrieval, 0 global scan of every host, >0 hostId

- **namingConventionForFiles:** naming convention for data files: long (complete naming, may lead to too long file names), short (abbreviated naming, may also lead to too long file names), hash (hash based naming, compact but may lead to naming clashes)
- **export:** export a replicate of the experiment with elements designated as a string with tokens separated by white spaces (addition of 'z' outputs a zip file): 'c' (code), 'd' (dependencies), 'i' (input), 1 (output data of step 1), 1d (stored data of step 1 only), 1o (observations of step 1 only)
- **sync:** sync data across machines specified as a cell array with: elements (as in bundle), optional host as numeric id (1.2 means the second host of the first group of machines, default to 2), direction 'd' (download from server to local host, default), 'u' (upload from host to server). Provided with the elements as string, the command default to server 2 in download mode.
- **clean:** clean explains directories and project data repositories. As string id : t (expLanes temporary directory), b (experiment backup directory), k (all steps directories while keeping reachable settings), 1 (output data of step 1), 1d (remove stored data of step 1 only), 1o (observations of step 1 only). As numeric id (clean the corresponding step directory, 0 means all directories)
- **branchStep:** specify at which step the branching (if any) is effective (default 0)

### *Exposition*

- **display:** step for which to display observations: -1 none, 0 last processed step if any, >0 specific step
- **expose:** specify default display of observations (>: prompt, t: table, p:plot)
- **tableDigitPrecision:** mantissa precision for the display of observations
- **displayFontSize:** font size for the display of observations in figures

### *Report*

- **report:** generate report: combination of r (run report), c (LaTeX compilation), and d (debug output)
- **reportName:** name of report to compile: empty (default), if containing the word <slides> or <Slides>, beamer presentation mode is used
- **latexDocumentClass:** style of LaTeX document (warning: this parameter is taken into account only at the creation of the LaTeX report file)

- **pdfViewer**: path to the pdf viewer (if left empty, expLanes will try to locate it automatically)
- **significanceThreshold**: threshold for statistical significance testing
- **showFactorsInReport**: show the factors graph in the report: 0 no display, 1 compact display, 2 also show stored data, 3 also show observations, 4 also show stored data and observations (if negative only generate the figure for latter inclusion)
- **factorDisplayStyle**: style of the factor graph: 0 no propagation of factors, 1 propagation of factors, 2 propagation of factors with an "all steps" node
- **figureCopyPath**: path where saved figures are copied outside the expLanes project
- **tableCopyPath**: path where saved tables are copied outside the expLanes project

### *Miscellaneous*

- **useExpCodeSmtplib**: usage of default expcode smtp to send emails. If set to 0, assume availability of local smtp server with complete credentials
- **sendMail**: send email: -2 server mode, email at start and end, -1 server mode, email at end only, 0, no email ever, 1 one email at end, 2 email at start and end

## Exposition of observations

To generate Tables and Figures that can be saved or embedded in your reports, the main function to use is:

```
config = expExpose(config, '<typeOfExposition>', ...
<parameters>);
```

The type of exposition specifies the rendering style. Several predefined styles are available and custom ones can easily be built. Parameters customize the rendering and are defined as string / value pairs.

**Important:** as expLanes does not use global variables, the state of the experiment is entirely stored in the `config` variable. Therefore, this variable must be propagated through each of the exposition calls.

### *Available expositions*

Available expositions are:

- **>**: redirect the observations to the command prompt
- **l**: generate a  $\text{\LaTeX}$  table

- **t**: Matlab table
- **b**: bar plot
- **p**: line plot
- **s**: scatter plot with up to 4 observations respectively mapped to the x-axis, the y-axis, the point size and the color
- **x**: box plot
- **a**: anova display
- **i**: image display

### *Custom exposition*

Should you require a custom way to display observation, you can do so by calling:

```
config = expExpose('custom');
```

At the first call, you are prompted:

```
Unable to find the function exposeCustom in your path.
```

This function is needed to display the observations with the Custom type of display.

```
Do you want to create it ? Y/N [Y]:
```

After acknowledgment, a skeleton script is generated in the codePath. This function has the following signature:

```
config = exposeCustom(config, data, p)
```

where

- **config**: contains the expLanes configuration state
- **data**: contains the observations as a struct array
- **p**: contains the exposition parameters

You can then fill the body of this function to suit your needs.

### *Parameters*

Parameters customize the rendering and are defined as string / value pairs.

- **addSpecification**: add display specification to the plot directive as ('parameter' / value) pairs; value can be a cell array, in this case the cell array is split across plot items
- **addSettingSpecification**: add display specification to the plot directive relative to specific settings as ('parameter' / value) pairs; value has to be a cell array which is split across plot items which correspond to the settings
- **caption**: caption of display as string; symbol "+" gets replaced by a compact description of the setting

- **color**: color of line; 1: default set of colors, 0: black, 'r', ...: user defined set of colors
- **compactLabels**: shorten labels by removing common substrings (default 0)
- **data**: specify data to be stored (default empty)
- **design**: use a canonical experimental design: assume cell array of [factors as numericIds], number of values per factor (0 complete set of values), type of plan as string 'f' (factorial) or 'o' (one factor at a time), [ optionally a seed as a vector of index of factor]. Shortcuts: numeric id (equivalent to [], <numericId>, 'f'), or 'one' design (equivalent to [], 2, 'o'), or 'star' design (equivalent to [], 0, 'o')
- **expand**: name or index of the factor to expand
- **fontSize**: set the font size of  $\LaTeX$  tables (default 'normal')
- **highlight**: highlight settings that are not significantly different from the best performing setting; -1: no variance, 0: variance for all observations, [1, 3]: variance for the first and third observations
- **highlightStyle**: type of highlighting; 'best': highlight best and equivalent (default), 'Best': highlight only the best, 'better': highlight the best if significantly better than the others, 'Better': highlight only the best if significantly better than the others
- **highlightColor**: use color to show highlights (default 1), -1 do not use \* to display the best performing one
- **integrate**: factor(s) to integrate by summation
- **label**:  $\LaTeX$  label of display as string (equal to the name if left empty)
- **legendLocation**: location of the legend (default 'BestOutSide', doc legend for more options)
- **marker**: specification of markers for line plot; 1: Matlab default set of markers (default), 0: no markers, "", ...: user defined cell array of markers
- **mask**: selection of the settings to be displayed
- **mergeDisplay**: concatenate current display with the previous one; "": no merge (default), 'h': horizontal concatenation, 'v': vertical concatenation
- **multipage**: activate the multipage setting for the  $\LaTeX$  table
- **name**: name of exported file as string; symbol "+" gets replaced by a compact description of the setting
- **number**: add a line number for each setting in tables



- **noFactor**: remove setting factors
- **noObservation**: remove observations
- **obs**: name(s) or index(es) of the observations to display
- **orderFactor**: numeric array specifying the ordering of the factors
- **orderSetting**: numeric array specifying the ordering of the settings
- **orientation**: orientation of display; 'v': vertical (default), 'h': horizontal, 'i': as second letter invert the table for prompt and latex display
- **percent**: display observations in percent; selector is the same as 'highlight'
- **plotCommand**: set of commands executed after the plotting directives as a cell array of commands
- **plotAxisProperties**: set of command setting the current axis after plotting directives as a cell array of couple property / value (see axis properties in Matlab documentation)
- **precision**: mantissa precision of data as numeric value or array; -1: take value of `config` parameter `tableDigitPrecision` (default), 0: no mantissa
- **put**: specify display output; 0: output to command prompt, 1: output to figure, 2: output to  $\text{\LaTeX}$
- **report**: include the current exposition in the report; 0: do not include (default for Matlab tables), 1: include (default for other displays)
- **rotateAxis**: rotate X axis labels (in degrees)
- **show**: display; 'data': actual observations (default), 'rank': ranking among settings, 'best': select best approaches, 'Best': select the significantly best approach (if any)
- **save**: save the display; 0: no saving, 1: save to a file with the name of the display as filename 'name': save to file 'name'
- **shortObservations**: compact observation names
- **shortFactors**: compact factor names
- **showMissingSetting**: show missing settings (default 0)
- **sort**: sort settings according to the specified set of observations if positive or to the specified set of factors if negative
- **step**: show observations for the specified processing step (as name or index). Defaults to the last step

- **title:** title of the display, as string, symbol “+” gets replaced by a compact description of the setting
- **total:** display average or summed values for observations; ‘v’, ‘V’: vertical with (v) or without (V) display of settings, ‘h’, ‘H’: horizontal with (h) or without (H) display of settings
- **variance:** display variance; selector is the same as ‘highlight’,
- **visible:** show the figure (default 1 except in save mode)

## Publishing your experiment

Publishing your experiment together with your research paper has two major benefits: not only is it the best way to assert beyond any possible doubt the trustworthiness of your reporting, it can also have a significant impact on the pace at which your research community progresses, as it allows other researchers to replicate exactly your experiments in an efficient way.

Indeed, unlike our introductory example, any real-life processing chain has many bells and whistles that needs to be tuned correctly in order to perform as described in your article. As it surely took you time and great care to reach that tuning, it is very unlikely that another researcher, as skillful as he may be, will dedicate the time needed to replicate it.

expLanes makes it possible to export a self-contained clone of your experiment, which can then be run on any compatible system, with this simple command:

```
geometricShape( 'export', 'c d', 'exportPath', ...
    <pathToExportLocation> );
```

### *Publishing input data*

Is the data used in your experiment

- freely available ?
- licensed with creative commons licenses ?
- your property and you have the ability to distribute it freely ?

If it falls into any of those categories, then your experiment can be replicated entirely which is, by far, the best case. If the data is not yet available online at the time of publishing, a convenient way to make data available online is the Archive repository<sup>4</sup>. Please add the needed documentation on how to retrieve the data online in the `README.txt` file of the experiment.

<sup>4</sup> <https://archive.org>

If the data is of small size, it can be conveniently embedded into the published version of the experiment with the ‘i’ token in the export command:

```
geometricShape( 'export', 'c d i', 'exportPath', ...
    <pathToExportLocation> );
```

*Publishing expLanes processing data*

If you do not wish to publish input data, you can still help others members of the community quickly gain information about your work by publishing expLanes processing data.

For example, if the first step of the experiment is to compute features, and you can ensure that the processing is not invertible, you may wish to publish the output of this first step, as it will allow others to replicate the following steps of the experiment:

```
geometricShape('export', 'c d 1d', 'exportPath', ...
    <pathToExportLocation>);
```

In that example, '1d' stands for the output data of step 1.

If you do not wish to do that, you can still export the observations of the processing steps in order to provide performance data that is not available in your research paper due to page limit constraints:

```
geometricShape('export', 'c d 1o', 'exportPath', ...
    <pathToExportLocation>);
```

Here, '1o' stands for the observations of step 1.

*Ensuring replicability**Procedure*

Thus, the procedure to safely publish replicable code using expLanes is:

## 1. Fill replicability information

- (a) start a Matlab session (if possible with another blank user account)
- (b) type `restoredefaultpath`, in order to remove any implicit dependencies
- (c) run the command you wish other users to execute to replicate your experiment:

```
geometricShape('do', 0, 'report', 'r', 'readMe', ...
    1);
```

at termination, the file `README.txt` will contain important informations about the software platform used, such as :

```
-----
Replication informations about the experiment ...
as of 30-Nov-2015
Command: geometricShape('do', 0, 'report', 'r');
Matlab version: 8.5.0.197613 (R2015a)
Loaded toolboxes:
- matlab
- signal_toolbox
- statistics_toolbox
-----
```

## 2. Create a duplicate:

```
geometricShape('export', 'c d i', 'exportPath', ...
    <pathToExportLocation>);
```

### 3. Test the duplicate

- (a) start a Matlab session (if possible with another blank user account)
- (b) go the export directory
- (c) optionally, decompress the zip archive
- (d) run the command available in the `README.txt` file.

## Advanced usage

### *Fork an experiment*

Experimental research most of the time involves the implementation of many solutions that will potentially improve things for a given step of the experiment. Finding the correct way to do it requires a lot of trials that may lead to code burden.

In order to keep a root project clean, and not to replicate the code and data and of this root project in a more experimental project, `explanes` allows you to create a fork experiment that accesses the data generated by the root experiment.

Let us assume that we want to create a new way of computing the volume of our geometric shapes. We can do so by creating a fork experiment:

```
expFork('forkGeometricShape', <pathToGeometricShape>, ...
2)
```

The last parameter specifies at which step the fork is done, in this case at step 2. The new experiment that replicate the structure of the root experiment in the directory `<pathToGeometricShape>` must keep the same factors up to the step before the fork step (in our case step 1).

The insertion of the new solution developed in the fork experiment into the root experiment is done manually.

### *Experiment design*

By default, `explanes` proceeds to the factorial exploration of every setting. The number of settings being the product of the cardinality of each set of modalities corresponding to the factors defined in the Factor file. For example, given 4 factors, each of 10 modalities, the number of settings to process is 10000, which may be impractical for any kind of large scale data processing.

One way is to sample the dataset, *i.e.* reduce the computation time needed for each setting. This can be done by using the dummy mode in the configuration of the experiment.

Another, more systematic, way is to prune out irrelevant sets of settings. To do so, one has to identify a minimal set of settings to process in order to quickly gain insights, given an set of observations of interest, about:

1. what are the most influential factors ?

## 2. are there any interaction between factors ?

Let us consider an experiment available in the `demo` directory called `designOfExperiments`. It has 4 factors (`f1`, `f2`, `f3`, ... `f4`), each with a set of modalities equal to 10. There is a unique observation of choice for this experiment (named `metric`) and the relationship between this observation and the factors is the following:

In this example, the answer to the 2 previous questions can easily be deduced. Factor `f4` does not need to be evaluated and factors `f1` and `f2` are correlated and thus have to be jointly considered.

Yet, in most systems, those facts cannot be easily deduced. Considering a reduced set of modalities allows the experiment to be tractable. This can be done by considering the `'design'` parameter. A factorial design with 2 values (minima and maxima) for each modality can be done using the following command

```
designOfExperiments('do', 1, 'design', 2);
```

which is equivalent to

```
designOfExperiments(('do', 1, 'design', [[]], 2, 'f'));
```

the empty array means consideration of each setting and 'f' stand for a factorial exploration.

```
designOfExperiments('expose', 'a', 'design', 2);
```

allows the user to study factors relationships through a multi-way factorial anova analysis, where the correlation between `f1` and `f2` and the uselessness of `f4` are shown. In order to identify the best setting, the remaining steps are:

```
designOfExperiments('do', 1, 'mask', {0 0 1 1});
```

The best performance is given for `f1=10` and `f2=10`. Then,

```
designOfExperiments('do', 1, 'mask', {10 10 0});
```

gives the best performance for `f3=10`. Using this method, the optimal setting (`10 10 0 *`) is identified in  $2^4 + 10^2 + 10 = 126$  computation of settings compared to  $10^4 = 10000$  for an exhaustive factorial exploration.

Alternative designs can be generated. The "one factor at a time" design allows you to explore the factors considering they are not correlated. In our example

```
designOfExperiments('design', 'one', 'do', 1);
```

is equivalent to

```
designOfExperiments('mask', {[[1 10] 1 1 1], {1 [1 ...  
10] 1 1}, {1 1 [1 10] 1}, {1 1 1 [1 10]}}], 'do', 1);
```

A "star" exploration can be done using

```
designOfExperiments('design', 'star', 'do', 1);
```

which is equivalent to

```
designOfExperiments('mask', {[0 1 1 1], {1 0 1 1}, ...  
{1 1 0 1}, {1 1 1 0}}, 'do', 1);
```

If available, an initial guess can be provided :

```
designOfExperiments('design', {[[]], 0, 'o', {5 4 3 ...  
2}}, 'do', 1);
```

	f1	f2	f3	f4
1 f1	0.00	0.00	0.37	0.37
2 f2	0.00	0.00	0.84	0.46
3 f3	0.37	0.84	0.00	0.54
4 f4	0.37	0.46	0.54	0.68

where the empty array stands for the exploration of each factor, `o` stands for the complete exploration of the modalities, `'o'` stands for "one factor at a time" exploration, and the last parameter is the (optional) seed. This last command is equivalent to:

```
designOfExperiments('mask', {{0 4 3 2}, {5 0 3 2}, ...
    {5 4 0 2}, {5 4 3 0}}, 'do', 1);
```

The design of experiments method can provide answers for other important questions. Barr & al. wrote a practical introduction on good practices in designing computational experiments <sup>5</sup>. Some books provide more in depth presentation <sup>6</sup>.

### *Handling compilation directives*

The `explanes` experiment may depends on mex files or external libraries that have to be compiled before the experiment is run.

If the code is static, and each processing host share the same computing environment, the compilation can be done once. If not, you may use the `expMex` function, placed in the `Init` file of the experiment: this function will, before any computation, proceed to the compilation if needed.

<sup>5</sup> Richard S. Barr, Bruce L. Golden, James P. Kelly, Mauricio G. C. Resende, and Jr. Stewart, William R. Designing and reporting on computational experiments with heuristic methods. *Journal of Heuristics*, 1(1):9–32, 1995

<sup>6</sup> Scott E. Maxwell, Harold D. Delaney, and Ken Kelley. *Designing Experiments and Analyzing Data: A Model Comparison Perspective, Second Edition*. Routledge, 2003; and George E. P. Box, J. Stuart Hunter, and William G. Hunter. *Statistics for Experimenters: Design, Innovation, and Discovery, 2nd Edition*. Wiley-Interscience, 2005

## Recommended readings

- **Getting it right:** R&D Methods for Science and Engineering, Peter Bock, Academic Press
- **The Visual Display of Quantitative Information** Edward R. Tufte
- **Warning Signs in Experimental Design and Interpretation** Peter Norvig.

<http://norvig.com/experiment-design.html>