

## MyESL Directives

No	Directives	Descriptions
<b>Model Options</b>		
1.	Position 1 <text file>	A text file containing full paths for all gene sequence alignments in FASTA format
2.	Position 2	A phylogenetic tree in newick format with clade ID
3.	--clade_list <text_file>	A text file containing the list of all clades for which clade models are built. Each line in the file contains a clade ID.
4.	--gen_clade_list <lower, upper>	This directive is used when the input phylogeny does not contain predefined clades. It automatically generates clade IDs for all internal nodes whose descendant species count falls within the specified <b>lower</b> and <b>upper</b> bounds. By default, both values are set to 0, which results in assigning clade IDs to <b>all internal nodes</b> in the phylogeny.
5.	--auto_name_length <int>	Specifies the number of characters to extract from species names when generating internal node labels for clade IDs. The default is 3, meaning the first three letters of species names will be used to construct clade identifiers.
6.	--clade_size_cutoff_lower <num>	This directive is used when no predefined clades are present in the input phylogeny. It enables users to automatically generate candidate clades for model building by specifying a lower and upper bound on clade size (both numeric). Clades containing a number of species within the specified size range will be selected, and ESL models will be built for each of them.
7.	--clade_size_cutoff_upper	
8.	--class_bal <method>	This directive addresses <b>class imbalance</b> in clade modeling—i.e., when the number of species inside the focal clade differs significantly from the number outside. Without

No	Directives	Descriptions
		<p>correction, such an imbalance can bias the model toward the larger class, reducing predictive accuracy.</p> <p>MyESL provides four strategies for class balancing:</p> <ul style="list-style-type: none"> <li>• <b>up</b>: Upsampling the smaller class to match the size of the larger class.</li> <li>• <b>down</b>: Downsampling the larger class to match the smaller class</li> <li>• <b>weighted</b>: Applying inverse class weights during model training</li> <li>• <b>phylo</b>: Phylogeny-aware balancing that selects training samples based on evolutionary relationships</li> </ul> <p>Use this option to ensure balanced and unbiased model training, especially when clade sizes are uneven.</p>
9.	<code>--classes &lt;text file&gt;</code>	<p>A tab-separated text file listing the <b>response values for all species</b> included in the clade model. Each line should contain a <b>species name</b> followed by its corresponding <b>binary response label</b>:</p> <ul style="list-style-type: none"> <li>• +1 for species <b>inside</b> the focal clade</li> <li>• -1 for species <b>outside</b> the clade</li> </ul> <p>This file is required when building clade-specific models.</p> <p>This file serves as the basis for building clade-specific models in the absence of a Newick-format phylogeny.</p>
10.	<code>--lambda1 &lt;float [0-1]&gt;</code>	<p>Specifies the site sparsity penalty parameter that controls the inclusion of features (e.g., sites) in the clade model. The value must be a float between 0 and 1, where higher values apply stronger penalization, resulting in a sparser model that selects fewer sites.</p>

No	Directives	Descriptions
11.	<code>--lambda2 &lt;float [0-1]&gt;</code>	Specifies the group sparsity penalty parameter that controls the inclusion of feature groups (e.g., genes) in the clade model. The value must be a float between 0 and 1, where higher values apply stronger penalization, resulting in a sparser model that selects fewer groups.
12.	<code>--no_group_penalty</code>	Enables monolevel Lasso analysis by disabling the group-level sparsity penalty. This directive must be used in combination with <code>--lambda1</code> , and ensures that sparsity is applied only at the individual feature level, not at the group (e.g., gene) level.
13.	<code>--lambda1_grid &lt;float, float, float [0-1]&gt;</code>	Specifies a range of feature-level sparsity parameters ( <code>lambda1</code> ) for building multiple clade models, one model for each parameter. This directive takes three float values: the starting value, the maximum value, and the step size. All values should be between 0 and 1. The parameter values increase linearly from the start to the maximum using the specified step, allowing exploration of different levels of feature sparsity.
14.	<code>--lambda1_grid &lt;float, float, float [0-1]&gt;</code>	Specifies a range of feature-level sparsity parameters ( <code>lambda1</code> ) for building multiple clade models, one model for each parameter. This directive takes three float values: the starting value, the maximum value, and the step size. All values should be between 0 and 1. The parameter values increase linearly from the start to the maximum using the specified step, allowing exploration of different levels of feature sparsity.

No	Directives	Descriptions
15.	<code>--grid_rmse_cutoff&lt;float&gt;</code>	This directive is to be used when lambda grids are used with multiple models. This helps to reduce the number of models to build and which models to use for aggregation.
16.	<code>--grid_acc_cutoff &lt;float [0 1]&gt;</code>	Sets the minimum training accuracy threshold for retaining models during grid search. Any model with accuracy below this cutoff will be discarded.
17.	<code>--grid_summary_only</code>	This directive retains only the summary of all models generated during the grid search and discards the results of individual models.
18.	<code>--min_group &lt;num&gt;</code>	Specifies the minimum number of groups that must be retained in each model during grid search. Models with fewer groups will be discarded.
19.	<code>--method &lt;method&gt;</code>	Specifies the type of loss function or regression to be used in model training. The default option is <code>logistic</code> . <ul style="list-style-type: none"> <li>• <code>logistic</code>: which performs sparse group lasso or lasso regression with logistic loss and is suitable for binary response variables (e.g., +1/-1).</li> <li>• <code>leastr</code>: applies sparse group lasso or lasso regression using least squares loss, and should be used when the response variable is continuous, ranging from <math>-\infty</math> to <math>+\infty</math>.</li> </ul>
20.	<code>--slep_opts &lt;text file&gt;</code>	A text file containing different options used for model optimization.
<b>Data Preprocessing</b>		
21.	<code>--include_singletons</code>	By default, MyESL removes all conserved sites and singleton variants (sites with only one differing sequence). Use this option to retain singleton sites during clade model construction if they are relevant to your

No	Directives	Descriptions
		analysis.
22.	<code>--bit_ct &lt;num&gt;</code>	Excludes sites where mutations are observed fewer than <num> times across all sequences. This helps filter out extremely rare variants that may be uninformative or noisy.
23.	<code>--data_type &lt;type&gt;</code>	<p>Specifies the type of input sequence data MyESL will process. Supported options are:</p> <ul style="list-style-type: none"> <li>• <code>nucleotide</code>: For DNA sequences</li> <li>• <code>protein</code>: For amino acid (AA) sequences</li> <li>• <code>molecular</code>: For datasets containing both DNA and protein sequences</li> <li>• <code>universal</code>: Allows any IUPAC characters, additional symbols, and numeric values (0–9), treating each unique character, including numbers, as distinct. This mode is case-sensitive, unlike the others, which are not.</li> </ul> <p>Choose the appropriate type based on the format and diversity of your input data.</p>
<b>Model Output</b>		
24.	<code>--output &lt;name&gt;</code>	An output directory name can be defined by users. This directory will be created in the current working directory. If no output directory name is defined by users, a directory named “output” will be created by default.
25.	<code>--gene_display_limit &lt;num&gt;</code>	Specifies the maximum number of gene groups (columns) to display in the model grid generated by MyESL. The default is 20, or fewer if the number of selected groups is less than 20.
26.	<code>--species_display_limit &lt;num&gt;</code>	Specifies the maximum number of taxa (columns) from the clade of interest to display in the model grid generated by MyESL. The default is 20, or fewer if the number of selected taxa is less than 20.
27.	<code>--gene_display_cutoff</code>	Defines the number of genes or groups

No	Directives	Descriptions
	<num>	(columns) to display in the model grid. The default is 20, or fewer if fewer than 20 groups are selected.
28.	--m_grid <row,col>	Sets the number of rows and columns in the model grid displayed by MyESL. The default is 20 for both dimensions, or fewer if the available rows or columns are less than 20.
29.	--stats_out <PGHS>	<p>Specifies the types of sparsity statistics to output from the model. Users can include any combination of the following:</p> <p>P: Site-level (position) sparsity scores  G: Gene or group sparsity scores  H: Hypothesis sparsity score  S: Species prediction scores</p>
<b>Miscellaneous</b>		
30.	--DrPhylo	Enables DrPhylo analysis for a specified clade, either defined by clade ID using a phylogeny in a NEWICK format or via a response file in text format. DrPhylo builds multiple sparse models across combinations of site and group sparsity parameters. By default, it explores values from 0.1 to 0.9 (in steps of 0.1), generating 81 models. Users can customize this range using the --lambda_grid directive. To reduce computation, DrPhylo skips models that include fewer than three genes. It outputs summary statistics, including PSS, GSS, and HSS scores, as well as a model grid. The default grid size is 20×20 (20 rows for species and 20 columns for genes), which can be adjusted using the --m_grid <row, col> directive.
31.	--subsets <int>	Enables subset-based analysis by dividing feature groups (e.g., genes) into multiple subsets. Models are trained on each subset, and informative groups are selected from each. A final model is then built using the combined set of selected groups. If

No	Directives	Descriptions
		--subsets 0, MyESL automatically determines the optimal number of subsets based on the system's computational capacity.