

# Imports

```
1 md"# Imports"
```

```
1 import Pkg, Revise; Pkg.activate(Base.current_project())
```



```
Activating project at `/DATA/cossio/SAM/2024/SamApp2024.jl`
```



```
1 import Makie, CairoMakie
```

```
1 import CSV, HDF5
```

```
1 import FASTX, Infernal
```

```
1 import SamApp2024
```

```
1 import RestrictedBoltzmannMachines as RBMs
```

```
1 import Rfam
```

```
1 import StatsBase, KernelDensity
```

```
1 using BioSequences: LongRNA
```

```
1 using DataFrames: DataFrame
```

```
1 using Distributions: Gamma, logpdf, pdf, Poisson
```

```
1 using LinearAlgebra: Diagonal, eigen
```

```
1 using Makie: @L_str
```

```
1 using NaNStatistics: nansum
```

```
1 using Random: bitrand
```

```
1 using Statistics: cor, mean
```

```
1 using StatsBase: countmap
```

# Plot

```
1 md"# Plot"
```

```
1 Rfam_cm = Infernal.cmfetch(Rfam.cm(), "RF00162");
```

```
RF00162_seed_stk =
```

```
▶(out = "/tmp/jl_KJi92rmN2l", stdout = "/tmp/jl_g3y9wLe96B", stderr = "/tmp/jl_4vCtudi1iN")
```

```
1 RF00162_seed_stk = Infernal.esl_afetch(Rfam.seed(), "RF00162")
```

```
1 RF00162_seed_match_cols = findall(≠('.'), SamApp2024.stockholm_ss(RF00162\_seed\_stk.out));
```

```
RF00162_seed_afa =
```

```
▶(out = "/tmp/jl_cVlIrYBgnex", stdout = "/tmp/jl_6NchYpGVWB", stderr = "/tmp/jl_GhKROx2XWB")
```

```
1 RF00162_seed_afa = Infernal.esl_reformat("AFA", RF00162\_seed\_stk.out;  
informat="STOCKHOLM") # WARNING: this has inserts marked as '-'
```

```
RF00162_seed_records =
```

```
▶[FASTX.FASTA.Record:                                     , FASTX.FASTA.Record:  
description: "AF027868.1/5245-5154"                        description: "AF269983.1/571-67
```

```
1 RF00162_seed_records = collect(FASTX.FASTA.Reader(open(RF00162\_seed\_afa.out)))
```

```
1 RF00162_seed_seqs_noinserts = LongRNA{4}.([FASTX.sequence(record)  
[RF00162\_seed\_match\_cols] for record in RF00162\_seed\_records]);
```

```
1 # trimmed (no inserts) aligned fasta
```

```
1 RF00162_hits_afa = Infernal.cmalign(Rfam\_cm.out, Rfam.fasta_file("RF00162");  
matchonly=true, outformat="AFA");
```

```
2 # these are already aligned and without inserts
```

```
1 RF00162_hits_sequences = LongRNA{4}.(FASTX.sequence.  
(FASTX.FASTA.Reader(open(RF00162\_hits\_afa.out))));
```

```
1 # emit sequences from Rfam CM model
```

```
2 Rfam_cm_emitted_sequences_afa = Infernal.cmemit(Rfam\_cm.out; N=5000, aligned=true,  
outformat="AFA");
```

```
▶[108nt RNA Sequence:                                     , 108nt RNA S  
CUUCUUAUUCAGAGGGACGGAGGGAGUGGCCUUGUAUAGU GCGAUACCCUUCUAG-----AUG--GCAAACGGAUAGGAG UUGCCAUUUGA
```

```
1 begin
```

```
2 Rfam_cm_emitted_sequences = FASTX.sequence.  
(FASTX.FASTA.Reader(open(Rfam\_cm\_emitted\_sequences\_afa.out)));
```

```
3 Rfam_cm_emitted_sequences = [filter(≠('.'), filter(!islowercase, seq)) for seq in  
Rfam_cm_emitted_sequences];
```

```
4 Rfam_cm_emitted_sequences = LongRNA{4}.(Rfam_cm_emitted_sequences);
```

```
5 end
```

► [108nt RNA Sequence: ACCCUAUC AAGAGUGGUUGAAGAGCUGGUCUUAUGAAAC...GCCAAGUCCGACAGGAUGGGAGUUCGGAAGAUGGGAG UACUUAUUCAC , 108nt RNA S

```
1 begin
2   # aligned hits, used to train a new noiseless CM model (in Stockholm format, without
   inserts!)
3   RF00162_hits_stk = Infernal.cmalalign(Rfam_cm.out, Rfam.fasta_file("RF00162"));
   matchonly=true);
4   # fit new CM model using full alignment (without inserts), and without entropic noise
5   Refined_cm = Infernal.cmbuild(RF00162_hits_stk.out; enone=true);
6
7   # emit sequences from Refined CM model
8   Refined_cm_emitted_sequences_afa = Infernal.cmemit(Refined_cm.cmout; N=5000,
   aligned=true, outformat="AFA");
9   Refined_cm_emitted_sequences = FASTX.sequence.
   (FASTX.FASTA.Reader(open(Refined_cm_emitted_sequences_afa.out)));
10  # remove inserts
11  Refined_cm_emitted_sequences = [filter(!=('.'), filter(!islowercase, seq)) for seq
   in Refined_cm_emitted_sequences];
12  @assert only(unique(length.(Refined_cm_emitted_sequences))) == 108
13  Refined_cm_emitted_sequences = LongRNA{4}.(Refined_cm_emitted_sequences);
14 end
```

```
1 # use saved RBM samples
2 sampled_v = SamApp2024.rbm2022samples(); # faster
```

► [119.98, 122.13, 117.93, 125.82, 126.6, 118.9, 118.78, 119.73, 114.7, 118.16, 121.37, 120.02, 12

```
1 begin
2   # Infernal scores of hits, using Rfam CM model
3   RF00162_hits_Rfam_cm_scores = Infernal.cmalalign_parse_sfile(
4     Infernal.cmalalign(
5       Rfam_cm.out,
6       Infernal.esl_reformat("FASTA", RF00162_hits_afa.out; informat="AFA").out;
7       glob=true, informat="FASTA"
8     ).sfile
9   ).bit_sc;
10
11  # Infernal scores of hits, using Refined CM model
12  RF00162_hits_Refined_cm_scores = Infernal.cmalalign_parse_sfile(
13    Infernal.cmalalign(
14      Refined_cm.cmout,
15      Infernal.esl_reformat("FASTA", RF00162_hits_afa.out; informat="AFA").out;
16      glob=true, informat="FASTA"
17    ).sfile
18  ).bit_sc;
19 end
```

► [105.91, 69.36, 103.03, 101.23, 107.89, 85.26, 60.42, 117.34, 102.88, 107.64, 98.22, 100.8, 110.

```
1 begin
2   # Infernal scores of Refined CM samples
3   _tmpfasta = tempname()
4   FASTX.FASTA.Writer(open(_tmpfasta, "w")) do writer
5     for (n, seq) in enumerate(Refined_cm_emitted_sequences)
6       ismissing(seq) && continue
7       write(writer, FASTX.FASTA.Record(string(n), filter(!=('-'), string(seq))))
8     end
9   end
10  # Infernal scores
11  Refined_cm_emitted_sequences_infernal_scores = Infernal.cmaligned_parse_sfile(
12    Infernal.cmaligned(Refined_cm.cmout, _tmpfasta; glob=true, informat="FASTA").sfile
13  ).bit_sc;
14
15
16  # Infernal scores of Rfam CM samples
17  _tmpfasta = tempname()
18  FASTX.FASTA.Writer(open(_tmpfasta, "w")) do writer
19    for (n, seq) in enumerate(Rfam_cm_emitted_sequences)
20      ismissing(seq) && continue
21      write(writer, FASTX.FASTA.Record(string(n), filter(!=('-'), string(seq))))
22    end
23  end
24  # Infernal scores
25  Rfam_cm_emitted_sequences_infernal_scores = Infernal.cmaligned_parse_sfile(
26    Infernal.cmaligned(Rfam_cm.out, _tmpfasta; glob=true, informat="FASTA").sfile
27  ).bit_sc;
28
29  # Infernal scores of RBM samples
30  _tmpfasta = tempname()
31  FASTX.FASTA.Writer(open(_tmpfasta, "w")) do writer
32    for (n, seq) in enumerate(SamApp2024.rnaseq(sampled_v))
33      @assert !ismissing(seq)
34      write(writer, FASTX.FASTA.Record(string(n), filter(!=('-'), string(seq))))
35    end
36  end
37
38  # Rfam CM
39  RBM_samples_Rfam_CM_infernal_scores = Infernal.cmaligned_parse_sfile(
40    Infernal.cmaligned(Rfam_cm.out, _tmpfasta; glob=true, informat="FASTA").sfile
41  ).bit_sc;
42
43  # Refined CM
44  RBM_samples_Refined_CM_infernal_scores = Infernal.cmaligned_parse_sfile(
45    Infernal.cmaligned(Refined_cm.cmout, _tmpfasta; glob=true, informat="FASTA").sfile
46  ).bit_sc;
47
48 end
```

```
1 # sites that have some non-zero fluctuations
2 # We need to separate frozen sites below because otherwise cor and eigen give NaN,
infinities, and fail
3 _variable_sites_flag = vec(all(0 .< mean(SamApp2024.onehot(RF00162_hits_sequences);
dims=3) .< 1; dims=1));
```

```
1 _variable_sites = findall(_variable_sites_flag);
```

```
1 RF00162_hits_var_sites_only = SamApp2024.onehot(RF00162_hits_sequences)[:,  
  _variable_sites, :];
```

```
1 RF00162_hits_cor = cor(reshape(RF00162_hits_var_sites_only, :,  
  size(RF00162_hits_var_sites_only, 3)); dims=2);
```

```
1 RF00162_hits_eig = eigen(RF00162_hits_cor);  
2  
3
```

```
1 # remap the variable sites eigenvectors back to the original consensus sequence numbering  
2 begin  
3   RF00162_hits_eigvec = zeros(5, 108, size(RF00162_hits_eig.vectors, 1))  
4   for n in 1:size(RF00162_hits_eig.vectors, 1)  
5     vec(view(RF00162_hits_eigvec, :, _variable_sites, n)) .=  
6     RF00162_hits_eig.vectors[:, n]  
7   end  
8 end
```

```
1 __proj_hits = reshape(SamApp2024.onehot(RF00162_hits_sequences), 5*108, :)' *  
  reshape(RF00162_hits_eigvec, 5*108, :);
```

```
1 __proj_rbm = reshape(sampled_v, 5*108, :)' * reshape(RF00162_hits_eigvec, 5*108, :);
```

```
1 __proj_refined_cm = reshape(SamApp2024.onehot(Refined_cm_emitted_sequences), 5*108, :)'  
  * reshape(RF00162_hits_eigvec, 5*108, :);
```

```
1 __proj_rfam_cm = reshape(SamApp2024.onehot(Rfam_cm_emitted_sequences), 5*108, :)' *  
  reshape(RF00162_hits_eigvec, 5*108, :);
```

```
1 # load SHAPE data  
2 shape_data_045 = SamApp2024.load_shapemapper_data_pierre_demux_20230920(; demux=true);
```

```
1 shape_data_rep0 = SamApp2024.select_conditions_20231002(shape_data_045,  
  filter(endswith("_rep0"), shape_data_045.conditions));
```

```
1 shape_data_rep45 = SamApp2024.select_conditions_20231002(shape_data_045,  
  filter(endswith("_rep45"), shape_data_045.conditions));
```

```
1 _idx_not_missing_seqs = findall(!ismissing, shape_data_rep0.aligned_sequences);
```

```
1 shape_sequences_onehot = SamApp2024.onehot(LongRNA{4}.  
  (shape_data_rep0.aligned_sequences[_idx_not_missing_seqs]));
```

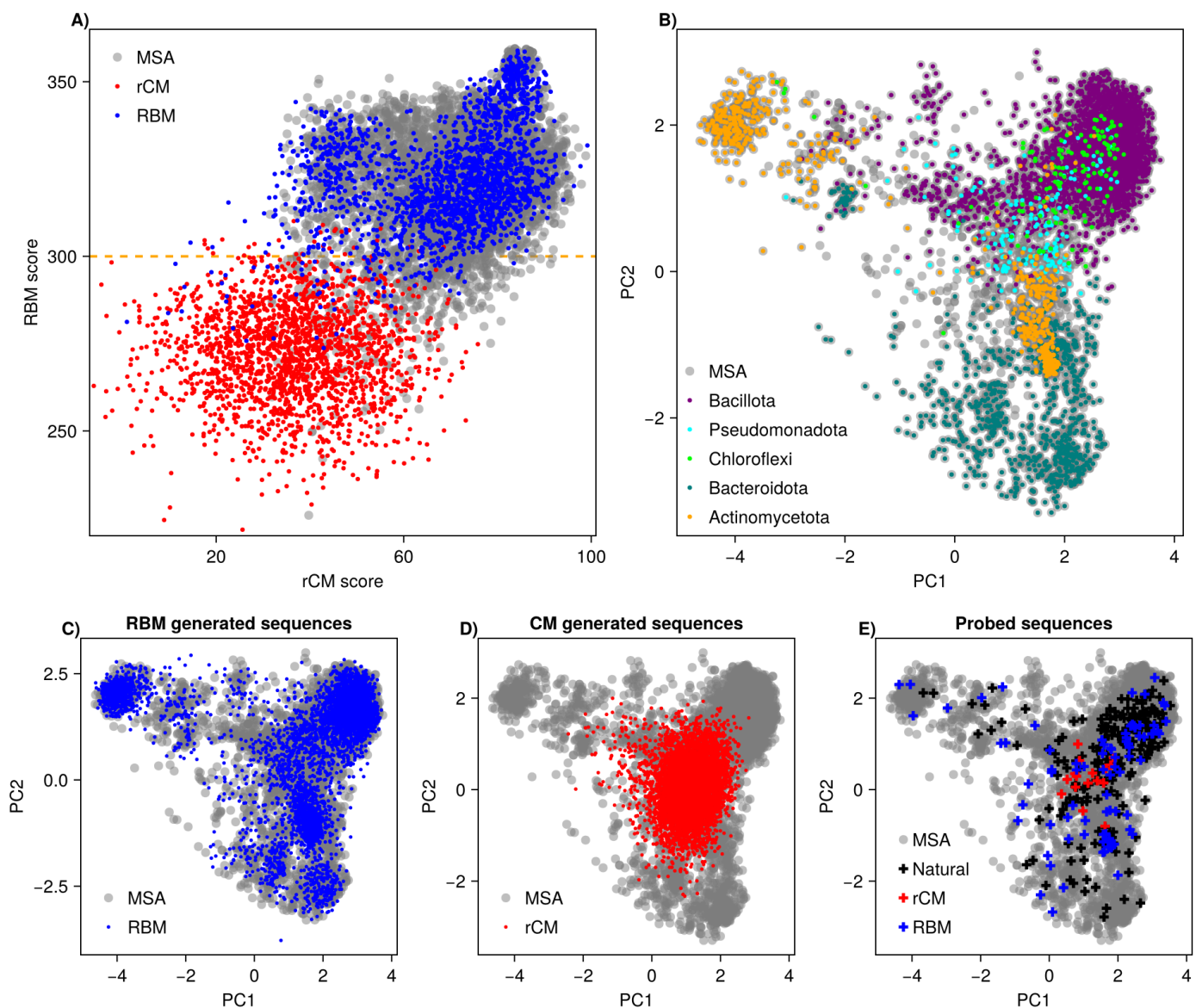
```
1 __proj_probed = reshape(shape_sequences_onehot, 5*108, :)' *  
  reshape(RF00162_hits_eigvec, 5*108, :);
```

```
1 _probed_origin = shape_data_rep0.aptamer_origin[_idx_not_missing_seqs];
```

►["Bacilli", "Bacilli", "Tissierellia", "Clostridia", "Bacilli", "Bacilli", "Bacilli", "Bacilli"]

```
1 begin
2   hits_tax_df = SamApp2024.rf00162_hits_taxonomy();
3   hits_tax_df.taxonomy_split = [ismissing(tax) ? missing : split(tax, "; ") for tax in
   hits_tax_df.taxonomy]
4   hits_tax_df.taxa_1 = [ismissing(tax) ? missing : length(tax) ≥ 1 ? tax[1] : missing
   for tax in hits_tax_df.taxonomy_split];
5   hits_tax_df.taxa_2 = [ismissing(tax) ? missing : length(tax) ≥ 2 ? tax[2] : missing
   for tax in hits_tax_df.taxonomy_split];
6   hits_tax_df.taxa_3 = [ismissing(tax) ? missing : length(tax) ≥ 3 ? tax[3] : missing
   for tax in hits_tax_df.taxonomy_split];
7 end
```

```
1 hits_tax_cnt = countmap(split(join(filter(!ismissing, filter(!ismissing,
   hits_tax_df.taxonomy)), "; "), "; "));
```



```

1 let fig = Makie.Figure()
2
3 ax = Makie.Axis(fig[1,1][1,1], xlabel="rCM score", ylabel="RBM score", width=400,
4 height=400, xticks=20:40:130, xgridvisible=false, ygridvisible=false)
5 Makie.hlines!(ax, 300, color=:orange, linestyle=:dash, linewidth=2)
6 #Makie.scatter!(ax, RF00162_hits_Refined_cm_scores, -
7 RBMs.free_energy(SamApp.rbm2022(), SamApp.onehot(RF00162_hits_sequences)),
8 label="Natural", color=:gray, 0.5), markersize=10)
9 Makie.scatter!(
10 ax, RF00162_hits_Rfam_cm_scores, -RBMs.free_energy(SamApp2024.rbm2022(),
11 SamApp2024.onehot(RF00162_hits_sequences)), label="MSA", color=:gray, 0.5),
12 markersize=10
13 )
14 Makie.scatter!(ax,
15 #Refined_cm_emitted_sequences_infernal_scores[1:2000],
16 Rfam_cm_emitted_sequences_infernal_scores[1:2000],
17 -RBMs.free_energy(SamApp2024.rbm2022(),
18 SamApp2024.onehot(Refined_cm_emitted_sequences))[1:2000],
19 label="rCM", color=:red, markersize=5
20 )
21 Makie.scatter!(ax,
22 #RBM_samples_Refined_CM_infernal_scores[1:2000],
23 RBM_samples_Rfam_CM_infernal_scores[1:2000],
24 -RBMs.free_energy(SamApp2024.rbm2022(), sampled_v)[1:2000],
25 label="RBM", color=:blue, markersize=5
26 )

```



```

21 Makie.xlims!(ax, -7, 101)
22 Makie.ylims!(ax, 220, 365)
23 Makie.axislegend(ax, position=:lt, framevisible=false)
24
25
26 _colors = [:purple, :cyan, :lime, :teal, :orange]
27 _c = 0
28
29 # Natural
30 ax = Makie.Axis(fig[1,1][1,2], xlabel="PC1", ylabel="PC2", width=400, height=400,
xgridvisible=false, ygridvisible=false) #title="Natural sequences")
31 Makie.scatter!(ax, __proj_hits[:, end], __proj_hits[:, end - 1], markersize=10,
label="MSA", color=:gray, 0.5))
32 for t = unique(hits_tax_df.taxa_2)
33     ismissing(t) && continue
34     hits_tax_cnt[t] > 100 || continue
35     _c += 1
36     Makie.scatter!(ax,
37         __proj_hits[replace(hits_tax_df.taxa_2 .== t, missing => false), end],
38         __proj_hits[replace(hits_tax_df.taxa_2 .== t, missing => false), end - 1],
39         markersize=5, label=t, color=_colors[_c])
40 end
41 Makie.axislegend(ax, position=(-0.05, -0.03), framevisible=false)
42
43 ax = Makie.Axis(fig[2,1][1,1], xlabel="PC1", ylabel="PC2", width=250, height=250,
xgridvisible=false, ygridvisible=false, title="RBM generated sequences")
44 Makie.scatter!(ax, __proj_hits[:, end], __proj_hits[:, end - 1], markersize=10,
label="MSA", color=:gray, 0.5))
45 Makie.scatter!(ax, __proj_rbm[:, end], __proj_rbm[:, end - 1], markersize=4,
label="RBM", color=:blue)
46 Makie.axislegend(ax, position=:lb, framevisible=false)
47
48 ax = Makie.Axis(fig[2,1][1,2], xlabel="PC1", ylabel="PC2", width=250, height=250,
xgridvisible=false, ygridvisible=false, title="CM generated sequences")
49 Makie.scatter!(ax, __proj_hits[:, end], __proj_hits[:, end - 1], markersize=10,
label="MSA", color=:gray, 0.5))
50 Makie.scatter!(ax, __proj_rfam_cm[:, end], __proj_rfam_cm[:, end - 1], markersize=4,
label="rCM", color=:red)
51 Makie.axislegend(ax, position=:lb, framevisible=false)
52
53 ax = Makie.Axis(fig[2,1][1,3], xlabel="PC1", ylabel="PC2", width=250, height=250,
xgridvisible=false, ygridvisible=false, title="Probed sequences")
54 Makie.scatter!(ax, __proj_hits[:, end], __proj_hits[:, end - 1], markersize=10,
color=:gray, 0.5), label="MSA")
55 Makie.scatter!(ax,
56     __proj_probed[( _probed_origin .== "RF00162_seed70") .| ( _probed_origin .==
"RF00162_full30"), end],
57     __proj_probed[( _probed_origin .== "RF00162_seed70") .| ( _probed_origin .==
"RF00162_full30"), end - 1],
58     markersize=10, color=:black, label="Natural", marker=:cross
59 )
60 Makie.scatter!(ax,
61     __proj_probed[_probed_origin .== "RF00162_syn_inf", end],
62     __proj_probed[_probed_origin .== "RF00162_syn_inf", end - 1],
63     markersize=10, color=:red, label="rCM", marker=:cross
64 )
65 Makie.scatter!(ax,
66     __proj_probed[_probed_origin .== "RF00162_syn_rbm", end],
67     __proj_probed[_probed_origin .== "RF00162_syn_rbm", end - 1],
68     markersize=10, color=:blue, label="RBM", marker=:cross

```



```
69 )
70 Makie.axislegend(ax, position=:lb, framevisible=false, patchlabelgap=-3)
71
72 Makie.Label(fig[1,1][1,1][1,1,Makie.TopLeft()], "A", font=:bold)
73 Makie.Label(fig[1,1][1,2][1,1,Makie.TopLeft()], "B", font=:bold)
74 Makie.Label(fig[2,1][1,1][1,1,Makie.TopLeft()], "C", font=:bold)
75 Makie.Label(fig[2,1][1,2][1,1,Makie.TopLeft()], "D", font=:bold)
76 Makie.Label(fig[2,1][1,3][1,1,Makie.TopLeft()], "E", font=:bold)
77
78 # fig[0,4] = Makie.Label(fig, "Natural sequences", font=:bold)
79 # fig[0,5] = Makie.Label(fig, "Generated sequences", font=:bold)
80
81 Makie.resize_to_layout!(fig)
82 #Makie.save("/workspaces/SamApp.jl/notebooks/2024-03-14 New paper
83 figures/Figures/PCA.pdf", fig)
```