

## Supplementary Methods

1. <i>Cultured human skeletal muscle cells</i> .....	2
2. <i>Mouse Muscle Injury</i> .....	2
3. <i>Human muscle biopsy processing</i> .....	3
a) Freezing of muscle tissues .....	3
b) RNA extraction protocol .....	3
c) Sample size requirements, RNA yield, and quality.....	7
4. <i>Differential gene expression</i> .....	9
5. <i>RNAseq-based classification</i> .....	12
a) Import data and normalize FPKM.....	12
b) Filter genes with low signal-to-noise ratio .....	13
c) Stratified cross-validation .....	13
d) Model training.....	13
e) Comparing the effect of including the gene-selection in the internal validation.....	15
f) Gene ranking based on the linear support vector machine classification.....	15
<i>References</i> .....	17

## 1. Cultured human skeletal muscle cells

Normal human skeletal muscle myoblasts (HSMM; Lonza) were cultured according to the manufacturer's protocol. When 80% confluent, the cultures were induced to differentiate into myotubes by replacing the growth medium with differentiation medium (DMEM, 2% horse serum, and L-glutamine). Two plates of cells were collected for RNA extraction at 7 separate time points: immediately before differentiation and then daily for 6 days.

## 2. Mouse Muscle Injury

Muscle injury and regeneration were induced in mice using cardiotoxin (CTX) as previously described.[1] Briefly, 6 week-old C57BL/6 mice were unilaterally injured by intramuscular injection of 0.1 mL of 10  $\mu$ M CTX into the tibialis anterior (TA) muscle. Injured TA muscles were harvested at days 3 (n=2), 5 (n=2), 7 (n=2), 10 (n=4), 14 (n=4), and 28 (n=3) post-injury. Contralateral (uninjured) TA muscles were also collected (n=9). Muscle tissue was snap-frozen and stored at -80 degrees Celsius.

### **3. Human muscle biopsy processing**

#### **a) Freezing of muscle tissues**

Open muscle biopsies were placed in an aluminum foil envelope. 2-methylbutane (isopentane) was pre-chilled using liquid nitrogen and the aluminum foil envelopes were submerged in the isopentane for 15 seconds. After this, the samples were placed in cryovials at -80° for long-term storage. Samples collected at other institutions were shipped in dry ice to the NIH Muscle Disease Unit.

#### **b) RNA extraction protocol**

##### **Required reagents**

1. TRIzol (ThermoFisher # 15596026)
2. Chloroform
3. 100% isopropanol
4. 75% ethanol
5. 1.4 ceramic bead homogenizing tubes (VWR # 10032-358)
6. RNase-free water
7. 1.5 mL low-binding tubes
8. GlycoBlue (ThermoFisher #AM9515)

**Step 0: Set-up reagents and equipment**

1. Spray down centrifuge, then set at 4°.
2. Set heat block at 55°C.
3. Set Precellys 24 homogenizer at 3x15" and 6500rpm.
4. Obtain dry ice and ice.
5. Thaw GlycoBlue.

**Step 1: Homogenization**

1. Get TRIzol from 4°.
2. Get homogenization tubes with beads.
3. Label tubes and set in dry ice.
4. Get samples from -80°C and place in dry ice.
5. Cut 1-2mm piece of the muscle biopsy with a #11 surgical blade.
6. Fill tubes in ice with 1mL TRIzol and immediately homogenize 3 x 15 minutes.

**Step 2: Phase separation**

1. Get chloroform, 100% isopropanol, 75% ethanol, and RNase-free water.
2. Wait for 5 minutes at room temperature.
3. Add 0.2 mL of chloroform.
4. Shake 15 seconds by hand.
5. Incubate 2-3 minutes at room temperature.
6. Centrifuge the sample at 12,000g for 15 minutes-4°C.
7. Label new 1.5mL tubes for all 4 samples.

8. Remove the aqueous phase samples, being careful not to touch interphase.
9. Place aqueous phases into new tubes.
10. Store interphase and organic phases -80°C.

### **Step 3: RNA Isolation**

1. Add 1.5uL of GlycoBlue, flicker and short spin.
2. Add 0.5 mL of 100% isopropanol.
3. Incubate at room temperature for 10 minutes.
4. Centrifuge the sample at 12,000g for 10 minutes at 4°C.

### **Step 4: RNA wash**

1. Remove supernatant (leave 0.5 mm).
2. Wash pellet with 1mL of 75% ethanol.
3. Vortex sample briefly.
4. Centrifuge the sample at 7,500g for 5 minutes at 4°C.
5. Discard the wash (with a pipette).
6. Dry RNA pellet for 5 minutes.

### **Step 5: RNA resuspension**

1. Resuspend RNA pellet in 40µL RNase-free water.
2. Incubate on heat block 10 minutes.
3. Measure the quantity of RNA with NanoDrop.

**Step 6: Evaluate RNA quality with 4200 TapeStation**

1. Allow High Sensitivity RNA Sample buffer (5067- 5580) to equilibrate at room temperature for 30 minutes.
2. Thaw High Sensitivity RNA ladder (5067- 5581) and total RNA samples on ice.
3. Launch the Agilent 4200 TapeStation Controller Software and select RNA assay mode under “settings”.
4. Flick the High Sensitivity RNA ScreenTape device (5067- 5579) and load it into the 4200 TapeStation instrument.
5. Place loading tips (5067- 5598) into the Agilent 4200 TapeStation instrument.
6. Vortex reagents and spin down before use.
7. Prepare diluted Ladder solution by adding 10  $\mu$ L RNase free water to the High Sensitivity RNA
8. Ladder vial and mix thoroughly. Pipette 1  $\mu$ L High Sensitivity RNA Sample Buffer and 2  $\mu$ L diluted High Sensitivity RNA Ladder at position A1 in a tube strip (401428).
9. For each sample, pipette 1  $\mu$ L High Sensitivity RNA Sample Buffer and 2  $\mu$ L RNA sample in a tube strip.
10. Cap tube strips with ladder or sample.
11. Mix liquids in sample and ladder vials using the IKA vortex at 2000 rpm for 1 min.
12. Spin down to position the sample and ladder at the bottom of the well plate and tube strip.

### 13. Samples and ladder denaturation:

- a. Heat samples and ladder to 72 °C (162 °F) for 3 min.
- b. Place samples and ladder on ice for 2 min.
- c. Spin down to position the samples and ladder at the bottom of the well plate and tube strip.

### 14. Sample Analysis:

- a. Load samples into the Agilent 4200 TapeStation instrument. Carefully remove caps of tube strips.
- b. Place the ladder in position A1 on tube strip holder in the 4200 TapeStation instrument.
- c. Select the required sample positions on the 4200 TapeStation Controller Software.
- d. Click “Start”.
- e. The Agilent TapeStation Analysis Software opens after the run and displays results.

### **c) Sample size requirements, RNA yield, and quality**

- The frozen muscle biopsy specimens were placed on glass pre-cooled with dry ice and a 1-2mm section of muscle was removed with a #11 surgical blade.

- An average of 11ug (SD 12ug) of RNA was recovered from each muscle biopsy specimen. 65ng of the samples were used to prepare the RNA library using the

NeoPrep™ system according to the TruSeq Stranded mRNA Library Prep protocol (Illumina) and sequenced using the Illumina HiSeq 2500 or 3000.

- The integrity of the RNA was verified using a standard quality metric denominated RNA integrity number (RIN) value, showing a median value of 7 (interquartile range [IQR] 5.9–7.4) for the muscle biopsy samples.

#### 4. Differential gene expression

We performed the differential expression between different subgroups using DESeq2 v.1.20.[2] DESeq2 performs an internal normalization where the geometric mean is calculated for each gene across all samples. The counts for a gene in each sample is then divided by this mean. The median of these ratios in a sample is the size factor for that sample. This procedure corrects for library size and RNA composition bias, which can arise for example when only a small number of genes are very highly expressed in one experiment condition but not in the other.

Additionally, DESeq2 automatically detects count outliers using Cooks' distance and removes these genes from the analysis. DESeq2 v.1.20 also performs independent filtering which maximizes the number of genes which will have a Benjamini and Hochberg-adjusted p-value less than a critical value set by default to 0.1; removing these genes with low counts improves the detection power by making the multiple testing adjustment of the p-values less severe. To speed up the computations we prefiltered genes with a total count across conditions below 10. Since these genes would have been excluded from the analysis afterward anyways, this did not influence the calculations at all.

DESeq2 uses shrinkage estimation for dispersions and fold changes. A dispersion value is estimated for each gene through a model fit procedure. Using these estimations, the package fits a negative binomial generalized linear model for each gene and uses the Wald test for significance testing. The Wald test P values

from the subset of genes that pass the independent filtering step are adjusted for multiple testing using the procedure of Benjamini and Hochberg.[3]

Loading the data, prefiltering and fitting the model was done with the following code:

```
1 library(DESeq2)
2 if (packageVersion("DESeq2") != "1.20.0") {
3   stop("DESeq2 version is not 1.20.0, please use version 1.20.0")
4 }
5
6 #Set working directory
7 project_folder <- "working_directory_path"
8 setwd(project_folder)
9
10 #Import count data and the sample covariates
11 countData <- as.matrix(read.csv("./anonymized_gene_counts.csv", row.names="gene_id"))
12 sample_covariates <- !(read.csv("./anonymized_gene_counts.csv", nrow=1, header = FALSE)[-1])
13 colnames(sample_covariates) <- "GROUP"
14
15 #Load data in DESeq2
16 dds <- DESeqDataSetFromMatrix(countData = countData, colData = sample_covariates, design= ~ GROUP)
17
18 #Prefiltering
19 dds <- dds[ rowSums(counts(dds)) >= 10, ]
20
21 #Fitting the model
22 deseq_fitted <- DESeq(dds, betaPrior=T)
```

To ensure the stability of the central tendency and dispersion values of each biological group between different sections of the study, the normalization process included the totality of the samples even if that specific comparison did not include some of those samples.

For example, the comparison between anti-Jo1 and normal biopsies was performed using the following code:

```
1 #Calculating the differential expression between anti-Jo1 and normal biopsies
2 deseq_results <- results(deseq_fitted, contrast=c("GROUP", "Jo1", "NT"))
3 write.csv(deseq_results[with(deseq_results, order(padj, pvalue)),], file=("./jo1_vs_nt.csv"))
```

We assigned equal weights to each autoantibody subgroups within DM and IMNM to avoid giving more importance to differentially expressed transcriptomic features of autoantibody subgroups with a higher number of biopsies at this stage of the analysis.

For example, to compare DM and IMNM we used the following code:

```
1 #Calculating the differential expression between DM and IMNM biopsies
2 resultsNames(deseq_fitted)
3 #'Intercept' 'GROUPHMGCR' 'GROUPIBM' 'GROUPJo1' 'GROUPMDA5' 'GROUPMi2' 'GROUPNT' 'GROUPNXP2' 'GROUPSRP'
  'GROUPTIF1'
4 deseq_results <- results(deseq_fitted, contrast=c(0,-1/2,0,0,1/4,1/4,0,1/4,-1/2,1/4))
5 write.csv(deseq_results[with(deseq_results, order(padj, pvalue)),], file=".dm_vs_imnm.csv")
```

## 5. RNAseq-based classification

### a) Import data and normalize FPKM

After importing the FPKM levels of all genes, we performed a logarithmic transformation and then normalized the data, maintaining the relative expression of the gene levels.

The rationale behind log-transforming the RNA expression values was to make variation similar across different orders of magnitude.[4]

As per the gene normalization, given that genes with low expression levels are more prone to technical bias in RNAseq,[5] the different genes were normalized so the relative expression of the gene levels would be preserved.

The following code was used to perform this step of the analysis:

```
1 import pandas as pd
2 import numpy as np
3
4 #Import data
5 df = pd.read_csv('anonymized_gene_fpkm.csv')
6
7 #Generate lists of samples
8 nt = df['GROUP'] == 'NT'
9 dm = df['GROUP'].isin(['Mi2', 'NXP2', 'TIF1', 'MDA5'])
10 asys = df['GROUP'] == 'Jo1'
11 immn = df['GROUP'].isin(['HMGCR', 'SRP'])
12 ibm = df['GROUP'] == 'IBM'
13 groups = {'nt':nt, 'dm':dm, 'as':asys, 'ibm':ibm, 'imnm':imnm}
14
15 df = df.set_index('GROUP')
16
17 #Log-transform FPKM values
18 df = np.log2(df+1)
19
20 #Normalize FPKM maintaining relative expression of gene levels
21 df = (df-np.mean(df.values))/np.std(df.values)
```

### b) Filter genes with low signal-to-noise ratio

To filter genes with low signal-to-noise ratio, we selected all the genes that were significantly different (with a cutoff of q-value <0.05) in each group compared to the rest. This was performed using the following lines of code:

```
1 #Filter genes with q-value>0.05
2 diff_expression_dfs = {}
3 sig_genes = []
4
5 for group in groups.keys():
6     diff_expression = pd.read_csv('./' + group + '_vs_all.csv')
7     diff_expression = diff_expression[diff_expression['padj'] < 0.05]
8     sig_genes = sig_genes + list(diff_expression['gene'])
9 sig_genes = list(set(sig_genes))
10 sig_genes = [gene.replace('-', '_').upper() for gene in sig_genes]
11 df = df[sig_genes]
12 df.shape
```

### c) Stratified cross-validation

We performed a stratified 3-fold-cross-validation of the samples included in the study. Stratification was used in order to ensure that there were enough samples in each cycle to build the models.

### d) Model training

With the objective of showing that the information contained in the RNAseq has classificatory value in myositis, we tested a set of classificatory machine learning models using the default parameters (for example, for AdaBoost we used the default algorithm AdaBoost-SAMME)[6]. The rationale behind performing this screening was that there were no comprehensive studies in the field to guide our selection of the best model for this application.

The cross-validation and training of the models were performed with the following lines of code:

```

1 import logging
2 from sklearn.model_selection import StratifiedKFold
3 from sklearn import metrics
4 from sklearn.svm import LinearSVC, SVC
5 from sklearn.ensemble import RandomForestClassifier, AdaBoostClassifier
6 from sklearn.neighbors import KNeighborsClassifier
7 from sklearn.naive_bayes import GaussianNB
8 from sklearn.tree import DecisionTreeClassifier
9 from sklearn.neural_network import MLPClassifier
10 from sklearn.gaussian_process import GaussianProcessClassifier
11 from sklearn.discriminant_analysis import QuadraticDiscriminantAnalysis
12
13 logging.basicConfig(level=logging.INFO, filename="models_log.log")
14
15 np.random.seed(1)
16
17 classifiers = {
18     "Linear SVM": LinearSVC(),
19     "RBF SVM": SVC(),
20     "Random Forest": RandomForestClassifier(),
21     "Nearest Neighbors": KNeighborsClassifier(),
22     "Gaussian Process": GaussianProcessClassifier(),
23     "Decision Tree": DecisionTreeClassifier(),
24     "Neural Net": MLPClassifier(),
25     "AdaBoost": AdaBoostClassifier(),
26     "Gaussian Naive Bayes": GaussianNB(),
27     "QDA": QuadraticDiscriminantAnalysis()
28 }
29
30 for model, clf in zip(classifiers.keys(), classifiers.values()):
31     results=[]
32
33     print(model)
34     logging.info(model)
35
36     for j in range(1000):
37         results_set = []
38
39         for i in groups.values():
40             #Diagnostic value
41             skf = StratifiedKFold(n_splits=3, shuffle=True)
42             skf.get_n_splits(df, i)
43
44             train_index, test_index = next(skf.split(df, nt))
45
46             try:
47                 clf.fit(df.iloc[train_index], i.iloc[train_index])
48                 y_pred=clf.predict(df.iloc[test_index])
49             except:
50                 pass
51
52             results_set.append(metrics.accuracy_score(i.iloc[test_index], y_pred))
53
54         results.append(np.array(results_set))
55         logging.info("{}0.2f".format(item) for item in results_set])
56
57     mean = np.mean(results, axis=0)
58     ci_lower = np.percentile(results, 2.5, axis=0)
59     ci_upper = np.percentile(results, 97.5, axis=0)
60
61     for n, i in enumerate(groups.keys()):
62         print(" {}: {2:0.1f} [{1:0.1f}-{3:0.1f}]" .format(i, ci_lower[n]*100, mean[n]*100, ci_upper[n]*100))

```

**e) Comparing the effect of including the gene-selection in the internal validation.**

We decided to use all the genes that were significantly different (with a cutoff of q-value  $<0.05$ ) in each group compared to the rest using all the samples. An alternative would have been to include the differentially expressed genes contained in the training set of each cycle. However, this approach was excessively computationally expensive. Nonetheless, to demonstrate the equivalency of these approaches, we modified our pipeline to train 100 cross-validation cycles using only the differentially expressed genes resulting from each training set. The performance of the models was equivalent using both methods (compare Table 2 with Supplementary Table 3).

**f) Gene ranking based on the linear support vector machine classification**

Once we determined that the linear SVM outperformed the rest of the models, we trained the whole dataset using this algorithm and then applied the recursive feature elimination technique [7] one gene at a time to sort the importance of the genes for each group of subjects.

The following lines of code were used to rank the genes for the different groups of subjects.

```
1 from sklearn.feature_selection import RFE
2
3 for name, dataset in zip(groups.keys(), groups.values()):
4     print("Top 10 variables " + name)
5
6     rfe = RFE(LinearSVC(), 1)
7     rfe = rfe.fit(df, dataset)
8
9     results = pd.DataFrame(sorted(zip(rfe.ranking_, df.columns)), columns = ['position', 'gene'])
10
11     results.to_excel('./svm_gene_importance_' + name + '.xlsx', index=False)
12
13     print(results.head())
```

## References

1. Mammen AL, Casciola-Rosen LA, Hall JC, Christopher-Stine L, Corse AM, Rosen A. Expression of the dermatomyositis autoantigen Mi-2 in regenerating muscle. *Arthritis Rheum*. 2009 Dec; 60(12):3784-3793.
2. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biol*. 2014; 15(12):550.
3. Benjamini Y, Drai D, Elmer G, Kafkafi N, Golani I. Controlling the false discovery rate in behavior genetics research. *Behav Brain Res*. 2001 Nov 1; 125(1-2):279-284.
4. Zwiener I, Frisch B, Binder H. Transforming RNA-Seq data to improve the performance of prognostic gene signatures. *PLoS One*. 2014; 9(1):e85150.
5. Bhargava V, Head SR, Ordoukhanian P, Mercola M, Subramaniam S. Technical variations in low-input RNA-seq methodologies. *Sci Rep*. 2014 Jan 14; 4:3678.
6. Zhu J, Zou H, Rosset S, Hastie T. Multi-class AdaBoost.
7. Guyon I, Weston J, Barnhill S, Vapnik V. Gene selection for cancer classification using support vector machines. *Mach Learn*. 2002; 46(1-3):389-422.