1

2

# Machine learning predicts translation initiation sites in

# neurologic diseases with expanded repeats

Short title: Machine learning, translation initiation, repeat expansion

disorders

7

8

9

10    Alec C. Gleason[1], Ghanashyam Ghadge[1,2], Jin Chen[3], Yoshifumi Sonobe[1,2], Raymond P. Roos[1,2*]

11

12    [1]University of Chicago, Chicago, Illinois, United States of America

13    [2]Department of Neurology, University of Chicago, Chicago, Illinois, United States of America

14    [3]Department of Pharmacology, University of Texas Southwestern Medical Center, Dallas, Texas,

15    United States of America

16

17    *Corresponding author

18    E-mail: rroos@neurology.bsd.uchicago.edu

# Abstract

A number of neurologic diseases, including a form of amyotrophic lateral sclerosis and others associated with expanded nucleotide repeats have an unconventional form of translation called repeat-associated non-AUG (RAN) translation. Repeat protein products accumulate and are hypothesized to contribute to disease pathogenesis. It has been speculated that the repeat regions in the RNA fold into secondary structures in a length-dependent manner, promoting RAN translation. Additionally, nucleotides that flank the repeat region, especially ones closest to the initiation site, are believed to enhance translation initiation. Recently, a machine learning model based on a large number of flanking nucleotides has been proposed for identifying translation initiation sites. However, most likely due to its extensive feature selection and limited training data, the model has diminished predictive power. Here, we overcome this limitation and increase prediction accuracy by a) capturing the effect of nucleotides most critical for translation initiation via feature reduction, b) implementing an alternative machine learning algorithm better suited for limited data, c) building comprehensive and balanced training data (via sampling *without* replacement) that includes previously unavailable sequences, and, d) splitting ATG and near-cognate translation initiation codon data to train two separate models. We also design a supplementary scoring system to provide an additional prognostic assessment of model predictions. The resultant models have high performance, with 85.00-87.79% accuracy exceeding that of the previously published model by >18%. The models presented here are then used to identify translation initiation sites in genes associated with a number of neurologic repeat expansion disorders. The results confirm a number of experimentally discovered sites of

40    translation initiation upstream of the expanded repeats and predict many sites that are not yet

41    established.

# Abbreviations

| RAN | Repeat-associated non-AUG |
| RLI | Repeat length-independent |
| KCS | Kozak consensus sequence |
| KSS | Kozak similarity score |
| AUROC | Area under receiver operating characteristic |
| ROC | Receiver operating characteristic |
| RFC | Random forest classifier |

# Introduction

More than 40 neurologic diseases are caused by expansions of repeat nucleotide sequences in causative genes. The repeats range from three nucleotides, such as 'CTG' associated with myotonic dystrophy Types I and II, to up to 12 nucleotides, such as 'CCCCGCCCCGCG', associated with progressive myoclonus epilepsy. Protein products translated from expanded repeat sequences tend to accumulate and aggregate, and have been proposed to contribute to disease [1-9]. Interestingly, in some cases, the repeats have been shown to be translated in all three reading frames from both the plus and minus strands of the RNA [10] by a process termed repeat-associated non-AUG (RAN) translation. It is believed that an affinity of translational machinery to folded regions of the RNA may underlie translation of the repeat sequences. Translation may occur from sequences in a repeat length-independent (RLI) mechanism. Regardless of repeat length, sequences may be ordered in such a way that they naturally increase the affinity of translational machinery to initiate at a particular codon. In such a process, translation may initiate not only within the repeat region, but also from sites upstream of the repeat sequences. In this case, repeat peptides will be produced if a stop codon is not encountered by the translational machinery before encountering the repeats. The large number of nucleotides that comprise and precede repeat sequences make the identification of RLI translation initiation sites challenging without proper laboratory evidence or computational methods.

A machine learning model called TITER has been proposed to predict all translation initiation sites in a given sequence. It addresses multiple limitations of the only other such model (to our best knowledge) [11, 12] and remains an important predictive tool. It appears, however, that the large feature selection of TITER and limited training data impair its predictive accuracy. The

65 predictive models described in our investigations have 85.00-87.79% accuracy that exceeds that

66 of TITER. Our models reduce the feature selection to capture the effect of ten critical nucleotides

67 that flank both sides of a putative translation initiation codon since a number of studies have

68 demonstrated a strong impact of nucleotides within this range on translation initiation [13-20].

69 We also introduce two models tailored for ATG or near-cognate codons because of their

70 differences in initiating translation [21, 22]. The models described here use an alternative

71 machine learning algorithm better suited for limited data [23]. We also present unbiased training

72 data through sampling techniques *without* replacement, using gene sequences that have been

73 unavailable to TITER. Finally, we generate a scoring metric to supplement model predictions.

74 The models confirm nearly all experimentally established translation initiation sites upstream of

75 repeats and, importantly, predict multiple sites that have not yet been investigated.

76

77

## Results

78

## Kozak similarity score algorithm

79

80 Before applying machine learning, we evaluated the performance of a more straightforward

81 algorithm that uses a limited number of nucleotides as predictors of translation initiation. This

82 algorithm was designed to predict the ability of a codon to initiate translation based on the

83 similarity of its surrounding sequence profile to the Kozak consensus sequence (KCS). The KCS

84 is a nucleotide motif, identified to most frequently border the canonical translation initiation

85 codon (ATG) and optimize translation initiation at the site. Although there exist slight variations,

6

86    this motif is typically accepted as the conserved pattern of the following underlined nucleotides

87    bordering the AUG codon: <u>CCRCC</u>AUG<u>G</u>. The nucleotide designated by R is a purine, most

88    typically adenine [13].

89    The sequence logo of the KCS (Fig 1) has been used to produce weighted scorings of identified

90    translation initiation codons and observe notable trends. The sequence logo illustrates conserved

91    nucleotides that tend to border ATG codons that initiate translation. The vertical length of each

92    letter in the sequence logo is related to the observed probability for a particular nucleotide to be

93    at a certain position, as well as the impact of the position on the efficiency of translation

94    initiation. It is formulated by the Shannon method [24].

95

96    **Fig 1. Schematic of the Kozak Similarity Score Algorithm.** Based on the sequences flanking an input

97    codon, the algorithm references the KCS Sequence Logo to assign the codon a score.

98

99    We designed a weighted scoring algorithm based on the KCS sequence logo and the ten bases

100   preceding and following a codon. Each nucleotide of the 23-base sequence has a value assigned

101   equal to the height of the nucleotide at its respective position, as illustrated in Fig 1. If a

102   nucleotide is not present in a position, it is assigned a value of zero. These values are then

103   summated, and the total divided by the maximal possible summated score (had each nucleotide

104   in the sequence been assigned the largest possible value for its position). This division serves to

105   make final values more feasible for interpretation. As opposed to the pre-normalized score range

106   of about 0 to 0.5990, scores derived from the normalization procedure more conveniently range

107    from 0 to 1. Overall, the final output score, referenced as Kozak similarity score (KSS), of a

108    codon is deduced by expression:

$$KSS(codon) \ = \ \frac{1}{KSS_{max}} \sum_{p=1}^{20} bits(nucleotide_p)$$

109    In this expression, $p$ denotes the position of a nucleotide bordering the codon. Values $p=1, 2, 3,$

110    $..., 10$ designate the positions of the ten nucleotides (from left to right) on the left side of the

111    codon, whereas values $p=11, 12, 13, ..., 20$ designate the positions of ten nucleotides (from left

112    to right) on the right side of the codon. Furthermore, *bits(nucleotide)* is the assigned height of a

113    particular nucleotide with reference to the KCS sequence logo (Fig 1). $KSS_{max}$ is the maximum

114    possible KSS that can be calculated for a codon.

115    We then used this algorithm on the sequences flanking known instances of ATG translation

116    initiation and produced a histogram distribution of the resulting scores (Fig 2). We created two

117    baselines to compare the scoring of ATG translation initiation codons against ATG codons that

118    do not initiate translation. For the first baseline, we ran the algorithm on one hundred thousand

119    'dummy' ATG codons that had completely randomized sequences without missing nucleotides

120    (a randomized adenine, cytosine, thymine, or guanine in every position flanking the codons) and

121    graphed the resulting score distribution. For the second, we ran the algorithm on a series of ATG

122    codons derived from the human genome that are believed not to initiate translation.

123

124    **Fig 2. Kozak Similarity Scores of ATG Translation Initiation Codons Against Baseline.**

125

126    As these histograms were generated from large datasets, they could more accurately serve as

127    representations of algorithm scoring for respective codon classifications: codons that initiate

128    translation, mixture of codons that initiate translation and do not initiate translation, and codons

129    that do not initiate translation, respectively.

130    Of note, the histogram in Fig 2 representing a randomized combination of codons that initiate

131    and do not initiate translation, is centered at about 0.59 for both the mean and median. In

132    contrast, in the histogram representing ATG codons that initiate translation, we observed a left-

133    skewed distribution, with mean and median scores of about 0.73 and 0.74, respectively. In the

134    histogram representing ATG codons expected not to initiate translation, we observed a slightly

135    right-skewed distribution, with mean and median scores of about 0.52 and 0.53, respectively.

136    Although exact sequences bordering near-cognate initiation codons have not been identified, as

137    has been carried out for the canonical ATG initiation codon (the KCS), current literature points

138    out similarities between the two sequences. For instance, in a bioinformatics study that analyzed

139    sequences bordering forty-five mammalian near-cognate initiation codons (including CUG,

140    GUG, UUG, AUA, and ACG),  a guanine or cytosine has been shown to frequent the -6 position

141    (6 bases upstream of the codon) [25]. As shown in Fig 1, a guanine or cytosine is also most

142    prevalent in the KCS at this position. The same study also noted the presence of a purine

143    (adenine or guanine) in the -3 position from the codon, which are the two most likely nucleotides

144    to occur in the same position of the KCS [25]. In a study of CUG near-cognate codons, those that

145    most frequently initiated translation had an adenine in the -3 position [26]. Although the

146    frequencies of adenine and guanine in the -3 position of the KCS are similar, analysis suggests

147    that adenine is more conserved. For example, if the nucleotide weightings in the KCS are

148    analyzed, adenine is conserved in about 47% of cases at the position versus that of guanine, with

149 about 37% conservation. Both the bioinformatics study as well as a publication analyzing peptide

150 translation from CUG-initiating mRNA constructs show enhanced translation when guanine is at

151 the +4 position (1 base downstream of the initiation codon) [18, 25]. In the KCS, guanine is most

152 conserved at the +4 position as well.

153 Because of these similarities, we decided to apply the algorithm to score known near-cognate

154 codons that have been shown to initiate translation (Fig 3). Interestingly, distributions of all

155 results are left-skewed, visibly differing from results derived from scoring of 'dummy' codons

156 with randomized flanking sequences, as well as codons expected not to initiate translation. In

157 particular, the distribution of scores for known CTG codons has mean and median of about 0.69.

158 The distribution of scores for known GTG codons has mean and median of about 0.69 and 0.70,

159 respectively. And the distribution of scores for known TTG codons has mean and median of

160 about 0.65. These results are an indication that the KSS of near-cognate codons can be used to

161 predict their ability to initiate translation.

162

163 **Fig 3. Kozak Similarity Scores of Near-Cognate Translation Initiation Codons Against Baseline.**

164

165 To use the KSS as a predictor of translation initiation ability, a threshold score has to first be

166 determined. In this way, an algorithm could classify codons with a score above the threshold as

167 initiating translation, and below it, not initiating translation. To find the best threshold, virtual

168 simulations were run using different score cutoffs to classify already known ATG initiation

169 codons and ATG codons expected not to initiate translation. Since there are at least 12,603 cases

170 of known ATG initiation codons in contrast to at least 34,097 ATG codons believed not to

171 initiate translation, the data were first balanced. In this way, the cut-off derived would not bias

172  classifications of codons in favor of not initiating translation. Next, all possible cutoff values

173  were set, ranging from 0.580 to 0.700 by increments of 0.001. This range was determined by

174  contrasting distributions in Fig 2. For each of these cutoff values, one thousand simulations were

175  run classifying the data of 12,603 known ATG translation initiation  codons on a randomized

176  subset containing 12,603 of the total 34,097 non-initiating ATGs. Errors were averaged for the

177  one thousand runs at each cutoff value. A cutoff of about 0.64 had the most minimized error.

178  When tested on data containing the 12,603 known ATG-initiating codons and randomized

179  12,603 instances of non-initiating ATGs, the average accuracy of the model was about 79.85%.

180  The area under receiver operating characteristic (AUROC) score from one of the thousand model

181  simulations (selected at random) was calculated to be 0.876. This score is a useful metric as it

182  indicates the model's discriminatory ability. In the model context, it would correctly assign a

183  greater prediction value for a codon to initiate translation if it indeed were a translation initiation

184  codon 87.6% of the time [27]. A random classifier has a score of 0.5, whereas a perfect classifier

185  has a score of 1.0 [28]. This score is calculated as the area under the ROC curve. This is a

186  graphical illustration of the model's ability to correctly categorize positives (true positive rate)

187  against decreased discrimination (increased false positive rate).

188  As carried out in the case of ATG, the cumulative data of the CTG, GTG, and TTG codons was

189  used to deduce a cutoff value for the algorithm's scoring of all near-cognate codons. To identify

190  the best cutoff for near-cognate codons, the same simulation process was used as was carried out

191  for ATG codons.  Using this simulation method, with balanced near-cognate codon data

192  consisting of equal numbers of positives (near-cognate initiation codons) and negatives (near-

193  cognate codons that do not initiate translation), the best cutoff of the algorithm classification was

194  about 0.61 for near-cognate codons. After a thousand simulations, the algorithm revealed an

195    average accuracy of about 75.60% for classifying near-cognate codons as initiating translation or

196    not initiating translation. The AUROC score calculated from one randomly selected simulation

197    was 0.835.

198

199    **Fig 4. Error Classifying ATG and Near-Cognate Codon Ability to Initiate Translation Using Kozak**

200    **Similarity Score.**

201

202    **Fig 5. ROC Curves of the ATG and Near-Cognate Kozak Similarity Score Classifiers.** The AUROC

203    score (area under the curve) of the ATG classifier is equal to 0.876. The AUROC score of the near-

204    cognate RFC is equal to 0.835.

205

206

## KSS as a reference for likelihood of translation initiation

207

208    In the previous section, the weighted scoring algorithm based on the KCS was used as a model to

209    classify whether codons could initiate translation. However, one could question whether the

210    scores of the weighting system could also be used as a metric. To investigate this issue, 12,603

211    instances of ATGs that initiate translation and 34,097 ATGs believed not to initiate translation

212    were compiled. One thousand balanced test datasets, containing the 12,603 positive ATG

213    instances along with randomly sampled negative ATG instances of the same number, were

214    gathered. The average proportion of codons that initiate translation with a KSS exceeding

215    particular values, across all test datasets was determined. These KSS thresholds ranged from zero

216    to one by increments of 0.02. The proportion of ATGs that initiate translation had a positive

217    correlation with the KSS. In other words, a greater proportion of ATGs would initiate translation

218 with an increased score. This score appeared useful since one could approximate the proportion

219 of ATG codons that initiate translation with equal KSSs to a particular codon encountered.

220 The same evaluation was conducted for near-cognate codons to deduce if there was a similar

221 trend. The procedures previously applied to the ATG data were used for the cumulative total of

222 2,413 instances of near-cognate codons that initiate translation, and 141,071 instances of near-

223 cognate codons believed not to initiate translation. There was a positive correlation between the

224 proportion of near-cognate codons that initiate translation and the KSS. In fact, the trend was

225 quite similar to that obtained for ATG data. The KSS was not limited as a metric for ATG

226 codons, but could be used to estimate the likelihood of a near-cognate codon to initiate

227 translation as well.

228 The results of the analysis for ATG and near-cognate codons is shown in the graph and table of

229 Fig 6.

230

231 **Fig 6. Proportion of ATG and Near-Cognate Codons that Initiate Translation with KSSs Above**

232 **Certain Values.** The graph and table were both generated to depict the same results, evaluated from

233 balanced data, i.e., an equal background proportion of positives and negatives.

234

235

## Random forest classifiers

237 A strong and practical approach for identifying translation initiation codons also includes the

238 application of a machine learning model. Machine learning models are powerful, as they can

239 analyze large amounts of complex data, determine patterns and codependences that are difficult

240 to process by a human, and learn from mistakes to improve over time [29]. Although biological

241 pathways are often sophisticated and produce remarkably diverse data, machine learning models

242 can provide direction for such processes that are not completely understood.

243 We decided to implement a random forest classifier (RFC). This machine learning algorithm

244 typically produces good results with partly missing data, bears little impact from outliers, and

245 mitigates overfitting. Furthermore, the RFC is a highly preferred model in contemporary

246 genomics [30]. The RFC is based on many decision trees, typically generated from large subsets

247 of data. As each decision tree may split data differently in the classification process, the

248 averaging of many such trees reduces variance and helps avoid overfitting. With an overfit

249 model, data inputs that vary slightly from trained data could have volatile classifications that are

250 not reliable. The RFC, which implements the averaging process, may produce greater accuracy

251 than any one decision tree alone [31].

252 Accordingly, an RFC was implemented as a separate algorithm to elucidate whether codons

253 initiate translation. To create such an algorithm, the feature variables of codons for the RFC to be

254 trained on were first assigned. For an ATG classifier, these variables designated the ten

255 nucleotides that preceded the codon, and ten that followed it. This range was chosen as studies

256 suggest that alterations of bases in some of these positions are highly impactful, and may define

257 whether a flanked codon is an "optimal, strong, [or] moderate" translation initiation site [13-20].

258 Although secondary structures can influence translation, which are dependent on a number of

259 nucleotides that may far exceed our incorporated range, successful identification of feature

260 patterns may require exceptionally large amounts of training data that are currently unavailable.

261 This is because the number of training samples required to differentiate data increases

262 exponentially as the number of attributes in a model increases [32]. Since five features are

263 needed to designate whether a nucleotide at each position, $n$, is either adenine, guanine, cytosine,

264    thymine, or missing, $5^n$ distinct data (enough to cover all possible data variations) may be

265    required for a model to best approximate the impact of each nucleotide, for every position that is

266    considered. By having our models trained on a narrowed scope of nucleotides known to

267    influence translation initiation, we sought to optimize predictive power with limited data. For a

268    near-cognate codon classifier, we included additional features to designate the nucleotide in the

269    first base position of the codons (i.e., the underlined: CTG, GTG, TTG). This is because the

270    nucleotide at this position may significantly impact translation initiation from these codons [21,

271    22].

272    Using the package, imbalanced-learn, in Python, we created the RFC models [33]. The ATG

273    RFC was trained using an imbalanced set of 12,603 ATG codons known to initiate translation

274    (positives), and 3,433 of 34,097 generated distinct ATG codons that are believed not to initiate

275    translation (negatives). The set of 3,433 negatives consisted of the total of 1,805 sequences that

276    were not missing nucleotides, and 1,628 (i.e., ten percent fewer) randomly sampled negatives of

277    the remaining 31,697 that were missing nucleotides. We left out five percent of the total 3,433

278    negatives used (172 ATGs that do not initiate translation), as well as the same number of

279    positives (172 ATG translation initiation codons) from the training data to constitute our test

280    dataset. In this way, accuracy would be based on unbiased data that was balanced with 344

281    combined cases of equally occurring positives and negatives.

282    The accuracy of the RFC model on the balanced 344 cases was 87.79%. In other words, the

283    algorithm correctly categorized 302 of the 344 ATGs, based on the sequences flanking each

284    codon. This accuracy is high in comparison to the 79.85% accuracy achieved using the KSS-

285    based classifier. We also calculated the area under receiver operating characteristic (AUROC)

286    score of the model to be 0.948, which is high as well. Increasing the parameter value designating

287     the total number of decision trees included in the RFC had no visible effect on model

288     performance. Other parameters were also best left unchanged for optimal predictions.

289     The same procedure was used to create an RFC for near-cognate codons as carried out for ATG

290     codons, using data available for all near-cognate codons. To prevent imbalanced data bias in the

291     accuracy measurement for the near-cognate RFC, data that was equally representative of all near-

292     cognate codons was set aside to form the test dataset. As the model was trained on CTG, GTG,

293     and TTG initiation codons, twenty positives and negatives were randomly isolated for each of

294     these codons prior to training. When run on this separated, balanced set of 120 data points, the

295     trained near-cognate RFC performed with 85.00% accuracy. The AUROC score of the near-

296     cognate classifier was calculated to be 0.938.

297

298     **Fig 7. ROC Curves of the ATG and Near-Cognate Random Forest Classifiers.** The AUROC score

299     (area under the curve) of the ATG RFC is equal to 0.948. The AUROC score of the near-cognate RFC is

300     equal to 0.938.

301

302

303     ## Analysis of the TITER neural network as a benchmark

304     To our knowledge, there exist only two other models for predicting both ATG and near-cognate

305     translation initiation codons. The latest is the TITER machine learning algorithm [11], which

306     addresses limitations of the first model. We analyzed TITER as a benchmark to compare it with

307     the performance of our presented models.

308    TITER is a deep learning-based framework that predicts whether a codon initiates translation

309    based on the product of two calculations, which is termed TISScore. One constituent is based on

310    the frequency of the codon of interest (e.g., ATG, CTG, GTG, etc.) in the dataset to initiate

311    translation. The second involves the averaging of calculated scores for a codon with flanking

312    sequences across thirty-two neural networks. A large number of neural networks was used as

313    part of a bootstrapping technique to account for training data imbalance.

314    Although TITER has a high AUROC score of 0.891 [11], ROC curves can present an "overly

315    optimistic" evaluation of a model's performance "if there is a large skew in the class

316    distribution" [27, 28]. This evaluation is based on the true positive and false positive rates of the

317    model – and an imbalance of positives and negatives may distort its calculation [34]. One

318    questions whether the test sample of the model is skewed as it consists of 767 positive and 9,914

319    negative samples in total [11]. Although the authors noted special procedures to account for the

320    data imbalance of the training dataset, it is not clear if such procedures were used for the test

321    dataset.

322    Since TITER was open-source, TITER's accuracy was averaged across a hundred balanced

323    subsets from its test dataset. Using all 767 positive samples, 767 negatives were randomly

324    sampled from the 9,914 total negatives, across the hundred runs to account for the data

325    imbalance. Through this technique, the unbiased average of the model accuracy was calculated

326    to be 66.94%. This was the accuracy achieved by the best cutoff, 0.5, of the TISScore for

327    classification. When run on the same sequences comprising the RFC test datasets (with

328    sequences extended to include the additional features TITER was trained with), TITER

329    demonstrated 62.21% and 58.33% accuracy for ATG and near-cognate codons, respectively.

330    These values were lower than the 75.60% and 79.85% accuracy achieved using the KSS scoring

331    system for ATG or near-cognate codons, or the 85.00% and 87.79% accuracy achieved using

332    RFC models. The fact that TITER was trained with less data than the RFC models presented here

333    could account for reduced predictive power. Specifically, it was generated using 9,776 positive

334    samples and 94,899 negatives compared to the total 15,016 positives and 175,168 negatives used

335    for the RFCs.

336    The performance of TITER may also be a result of the large number of features that this machine

337    learning model incorporated. Although contemporary research suggests a few bases that flank a

338    codon greatly influence translation initiation from the site [13-20], TITER analyzes a total of two

339    hundred bases that flank each codon. Compared to our approach of analyzing ten preceding and

340    proceeding nucleotides, TITER may implement up to $180*5 = 900$ additional features. The

341    expression '180*5' is used because any one base at the 180 extra positions is represented by five

342    features to designate whether the base is adenine, guanine, cytosine, tyrosine, or is missing.

343    Although the TITER publication mentions feature reduction in the hidden layer of the neural

344    networks, it is not clear how much feature reduction occurred and whether features with

345    significant correlations were inadvertently reduced. An excess of features may decrease

346    effectiveness in machine learning because the number of training samples required to

347    differentiate the data increases exponentially as the number of attributes in a model increases.

348    Thus, predictive power is lost. In fact, this phenomenon is termed the "curse of dimensionality"

349    in Data Science [32].

350    In addition to feature reduction, our implementation of the random forest classifier, which is

351    more robust to outliers and erroneous instances (especially when data is limited), creation of two

352    models to account for properties of different data types (i.e., ATG codons versus near-cognate

353 codons), and use of sampling *without* replacement which preserves natural variations found in

354 data (in place of bootstrapping) could explain our improved model performance.

355

356 **Fig 8. ROC Curves of All ATG and Near-Cognate Classifiers Derived from Same Test Data.** All

357 classifiers were run on the ATG and Near-cognate RFC test datasets, and their ROC curves were

358 superimposed. The AUROC scores of the ATG and near-cognate RFCs are 0.948 and 0.938, respectively.

359 The AUROC scores of the ATG and near-cognate KSS classifiers are 0.857 and 0.787 on these test

360 datasets. TITER's AUROC scores are 0.622 and 0.603 for ATG and near-cognate codons, respectively.

361

362

# Model selection and integration into software

363

364 Of the two types of models created, the RFCs appeared the best model to use for predicting

365 translation initiation sites. With accuracy determined from the balanced test dataset for ATGs at

366 87.79% and for near-cognate codons at 85.00%, their performance exceeds that of the

367 straightforward KSS-based classifiers. To our best knowledge, the RFCs also outperform all

368 other models designed for the same function, including TITER, which they exceed by more than

369 18% in accuracy. As a next step, we decided to use the RFCs to identify repeat-length-

370 independent (RLI) translation initiation associated with neurologic diseases.

371 To do this, the RFC models were implemented into software. Developed in Python, the program

372 could be used to evaluate a total sequence consisting of the upstream region, followed by ten

373 nucleotide sequence repeats to represent the repeat expansion. Ten sequence repeats may be

374 adequate to capture the repeat expansion effect on translation initiation from upstream codons as

375    well as codons within the repeat expansion itself because ten nucleotide sequence repeats are at

376    minimum thirty bases long, and the integrated model only uses the ten bases that flank each side

377    of a codon for analysis. Nucleotides within this range have been shown to strongly impact

378    translation initiation [13-20].

379    The model can scan through each codon in the sequence and return a prediction from the

380    implemented RFCs. If a codon encountered is 'ATG,' then the ATG RFC with 87.79% accuracy

381    predicts whether it initiates translation based on the ten sequences flanking each side of the

382    codon. Otherwise, if the codon encountered is a near-cognate codon, then the near-cognate RFC

383    with 85.00% accuracy predicts whether it initiates translation via the same procedure. Next, the

384    program virtually simulates translation from each predicted codon and filters out those instances

385    in which a stop codon (TAG, TGA, or TAA) is encountered upstream of the repeat expansion.

386    This feature was implemented to remove codons from consideration if their initiated translation

387    would not reach the repeat expansion and produce the pathogenic repeat proteins that are

388    associated with neurologic disease. Then, the program would determine the repeated nucleotide

389    sequence that would be translated from each predicted initiation codon, as well as the associated

390    translation product. Finally, the program outputs a visualization of the input sequence, with

391    predicted codons color-coded to distinguish the associated product translated.

392    In the figures that follow, nucleotides have a bold font to distinguish initiation codons that the

393    software models were trained on. These codons include canonical start codon ATG, and near-

394    cognate codons CTG, GTG, and TTG. Because the features of the three near-cognate codons

395    were used to extrapolate classifications of the other, less researched near-cognate codons (AAG,

396    AGG, ACG, ATC, ATT, and ATA), it is possible to incur false predictions for these less studied

397    instances. Thus, these six near-cognate codons are designated only with color-coding without

398    bolding to denote that they should be acknowledged with less confidence. If there is an overlap

399    between predicted initiation codons (i.e., one or two nucleotides overlap between predicted

400    codons), the color of the overlapped region is the same as that of the next predicted codon to

401    prevent confusion. The overlapped region may or may not be bolded depending on whether the

402    software was trained on this next codon. We also output the KSSs of each predicted codon to two

403    decimal points, as the score could be a useful metric to evaluate translation initiation likelihood.

404    This may be approximated through comparison of KSSs of a codon to the reference table and

405    graph (Fig 6).

406

407    **Fig 9. An Example of the Formatting Scheme in Software Output.** This example shows predicted

408    codons that are color-coded based on their reading frame: 'ATT,' 'TTG,' 'CTG,' 'AGG,' 'GTG,' and

409    'CTG.' Codons that the models were trained on show up with bold formatting. If there is an overlap

410    between predicted initiation codons (i.e., one or two nucleotides overlap between predicted codons), the

411    color of the overlapped region is the same as the color of the next predicted codon.

412

413

## Software ability to identify known RLI translation initiation sites

415    After the software was completed, its ability to distinguish RLI translation initiation sites was

416    analyzed. We first identified translation initiation codons upstream of repeats in the following

417    genes in which RAN translation is known to occur: *C9orf72* (associated with amyotrophic lateral

418    sclerosis and frontotemporal dementia), *FMR1* (associated with fragile X and fragile X-

419    associated tremor/ataxia syndrome), *DM1* (associated with myotonic dystrophy type 1), and

420     *HDL2* (associated with Huntington disease-like 2) genes. These examples were used as

421     references for software performance. It should be noted that translation initiation codons

422     identified for DM1 were obtained from an experiment that implemented a slightly modified

423     version of the conventional DM1 antisense strand. The strand had been experimentally modified

424     to determine whether changes in its sequence could induce translation initiation from particular

425     codons [35]. Next, the associated upstream regions and repeat expansion sequences for each

426     gene, as recorded in the National Center for Biotechnology Information database, were input into

427     the software. Predictions were generated in order to determine whether they corresponded to

428     experimentally confirmed translation initiation codons (Table 1).

**Table 1**. Previously identified RLI translation initiation sites from publications.

| Gene | Codon | Number of Bases Upstream of Repeat | Peptide Repeat Translated | Kozak Similarity Score |
|---|---|---|---|---|
| *C9orf72* (Sense) [4] | AGG | 1 | Poly-GR | 0.66 |
| | CTG | 24 | Poly-GA | 0.69 |
| *C9orf72* (Antisense) [4] | ATG | 194 | Poly-PG | 0.61 |
| *FMR1* (Sense) [36, 37] | GTG | 11 | Poly-G | 0.70 |
| | ACG | 35 | Poly-G | 0.80 |
| | ACG | 60 | Poly-G | 0.71 |
| *DM1* (Antisense) with slightly modified sequence [35] | ATC | 7 | Poly-A | 0.61 |
| | ATG | 17 | Poly-S | 0.66 |
| | ATT | 23 | Poly-S | 0.74 |
| *HDL2* (Antisense) [35] | ATC | 6 | Poly-Q | 0.74 |

429     Comparison between the predictions and experimentally identified translation initiation codons

430     demonstrated high performance of the software. In fact, all translation initiation sites previously

431    identified across existing publications were correctly identified by the RFCs with one exception:

432    ATC, which was found experimentally to initiate translation in the modified *DM1* antisense

433    strand seven bases upstream of the repeat [35]. However, the near-cognate RFC model

434    successfully predicted all other instances of translation initiation from less researched near-

435    cognate codons. This accuracy is surprising considering that the near-cognate RFC model was

436    only trained on instances of CTG, GTG and TTG translation. As there was insufficient data to

437    train the model on less used near-cognate codons (ATA, ATC, ATT, AGG, ACG, and AAG),

438    predictions for these codons were extrapolated based on recognized patterns from CTG, GTG,

439    and TTG examples. However, for the same reason that they were not included in model training,

440    near-cognate codons that are not CTG, GTG, or TTG should be acknowledged with less

441    confidence in predictions, out of concern they may be false positives.

442

## Predicted Translation Initiation Sites Associated with Neurologic Diseases

443

444

445    Experimentally identified translation initiation codons for *C9orf72*, *FMR1*, *DM1,* and *HDL2*

446    were confirmed by the model presented here (Table 1, Figs 10 and 11). As the software

447    performed well, it was then used to predict translation initiation codons associated with repeats

448    in neurologic diseases that have not been experimentally identified. The software was also used

449    to make predictions for translation initiation codons for other genes with repeats associated with

450    neurologic repeat diseases, *HTT*, and *DM2* (Fig 12). Predicted translation initiation codons with

451    relatively high KSSs were noted for all analyzed genes (Table 2). In all cases, predicted

452    translation initiation sites are not shown if they have a downstream stop codon located in the

453    same reading frame before the repeat.

454

455    **Fig 10. Predicted Translation Initiation Codons for *C9orf72 and FMR1*.** Predicted codons that the

456    models were trained on show up with bold formatting. Numbers indicate the number of bases upstream of

457    the repeat.

458    * A predicted translation initiation codon overlaps with the repeat (AGG, located 1 base upstream).

459    **Fig 11. Predicted Translation Initiation Codons for *DM1 and HDL2*.** Predicted codons that the models

460    were trained on show up with bold formatting. Numbers indicate the number of bases upstream of the

461    repeat.

462    * Every CTG within the repeat is predicted to possibly initiate translation.

463    † Every CTG within repeat, aside from the first one, is predicted to possibly initiate translation.

464    **Fig 12. Predicted Translation Initiation Codons for *HTT and DM2*.** Predicted codons that the models

465    were trained on show up with bold formatting. Numbers indicate the number of bases upstream of the

466    repeat.

467    * Every CTG within the repeat, aside from the first one, is predicted to possibly initiate translation.

468    † Two predicted translation initiation codons are within repeat.

469

470

471

472

473

**Table 2**. Translation Initiation Codons with High Kozak Similarity Scores per Translated Polypeptide Repeat*

| Codon | Number of Bases Upstream of Repeat | Kozak Similarity Score | Translated Polypeptide Repeat |
|---|---|---|---|
| **C9orf72 (Sense)** | | | |
| CTG | 24 | 0.66 | Poly-GA |
| AGG | 1 | 0.69 | Poly-GR |
| **C9orf72 (Antisense)** | | | |
| **ATG†** | 113 | 0.75 | Poly-PG |
| AAG | 350 | 0.84 | Poly-PG |
| ACG | 3 | 0.79 | Poly-PR |
| AAG | 288 | 0.73 | Poly-PR |
| AAG | 384 | 0.77 | Poly-PR |
| **FMR1 (Sense)** | | | |
| AGG | 18 | 0.83 | Poly-R |
| ACG | 60 | 0.71 | Poly-R |
| ACG | 35 | 0.79 | Poly-G |
| **GTG** | 38 | 0.76 | Poly-G |
| AAG | 332 | 0.83 | Poly-G |
| **FMR1 (Antisense)** | | | |
| AGG | 28 | 0.71 | Poly-A |
| **GTG** | 26 | 0.73 | Poly-R |
| **CTG** | 56 | 0.70 | Poly-R |
| ATT | 105 | 0.81 | Poly-P |
| AAG | 156 | 0.78 | Poly-P |
| AAG | 177 | 0.85 | Poly-P |
| **CTG** | 195 | 0.74 | Poly-P |
| AGG | 207 | 0.84 | Poly-P |
| ATC | 252 | 0.80 | Poly-P |
| AGG | 318 | 0.74 | Poly-P |
| **DM1 (Sense)** | | | |
| AAG | 23 | 0.62 | Poly-C |
| AGG | 61 | 0.77 | Poly-A |
| **CTG** | -1 | 0.67 | Poly-L |
| **DM1 (Antisense)** | | | |
| **CTG** | 34 | 0.87 | Poly-A |
| AGG | 169 | 0.85 | Poly-A |
| ATC | 193 | 0.81 | Poly-A |
| ACG | 98 | 0.86 | Poly-S |
| **HDL2 (Sense)** | | | |
| ATC | 72 | 0.71 | Poly-L |
| ATC | 68 | 0.52 | Poly-C |
| AGG | 10 | 0.84 | Poly-A |
| **HDL2 (Antisense)** | | | |
| ATC | 6 | 0.74 | Poly-Q |
| AAG | 27 | 0.80 | Poly-Q |
| ATT | 261 | 0.81 | Poly-Q |
| **GTG** | 372 | 0.83 | Poly-Q |

| | | | |
|---|---|---|---|
| **GTG** | 378 | 0.71 | Poly-Q |
| **CTG** | 122 | 0.68 | Poly-S |
| ATC | 67 | 0.69 | Poly-A |
| *HTT* **(Sense)** | | | |
| AAG | 27 | 0.76 | Poly-Q |
| **CTG** | 33 | 0.72 | Poly-Q |
| **CTG** | 42 | 0.87 | Poly-Q |
| **ATG** | 51 | 0.89 | Poly-Q |
| AAG | 210 | 0.72 | Poly-Q |
| **CTG** | 348 | 0.74 | Poly-Q |
| ACG | 187 | 0.75 | Poly-A |
| **GTG** | 202 | 0.85 | Poly-A |
| *HTT* **(Antisense)** | | | |
| ATC | 213 | 0.76 | Poly-L |
| AGG | 225 | 0.70 | Poly-L |
| AAG | 330 | 0.73 | Poly-L |
| **CTG** | 342 | 0.70 | Poly-L |
| AGG | 369 | 0.76 | Poly-L |
| **GTG** | 13 | 0.84 | Poly-A |
| **GTG** | 118 | 0.72 | Poly-A |
| **CTG** | 199 | 0.81 | Poly-A |
| **CTG** | 229 | 0.71 | Poly-A |
| **GTG** | 337 | 0.73 | Poly-A |
| *DM2* **(Sense)** | | | |
| **CTG** | 7 | 0.50 | Poly-CLPA |
| **CTG** | -5 | 0.61 | Poly-LPAC |
| ATT | 87 | 0.66 | Poly-PACL |
| *DM2* **(Antisense)** | | | |
| AGG | 7 | 0.72 | Poly-GRQA |
| **GTG** | 58 | 0.70 | Poly-GRQA |
| ATA | 88 | 0.75 | Poly-GRQA |
| AGG | 47 | 0.71 | Poly-RQAG |
| AGG | 113 | 0.74 | Poly-RQAG |
| AGG | 15 | 0.72 | Poly-QAGR |

*Predicted codons are displayed that have KSSs above 0.70. If no KSSs within a reading frame are above 0.70, then the codon with the highest KSS is presented – as in the case of the *C9orf72* sense strand.

† Bolded codons represent codons that the RFCs were trained on.

474

475    Results displayed in the figures and table indicate translation initiation sites for proteins

476    translated from the repeat. Of note, the average KSS of all upstream predicted codons is about

477    0.66. With reference to the table in Fig 6, approximately 80% of ATG and near-cognate codons

478    with a score above 0.65 are estimated to initiate translation from a background population of

479    equally occurring translation initiation codons (positives) and codons believed not to initiate

480    translation (negatives).

481    With respect to the *C9orf72* sense strand upstream from the repeat, the software predicts a codon

482    to initiate translation of poly-GA, and another to translate poly-GR. Both of these codons have

483    been confirmed through experimentation [4]. In the antisense strand, there are ten codons that

484    could initiate translation of poly-PR, and six predicted with respect to poly-PG. The ATG located

485    194 bases upstream of the repeat expansion has been confirmed [4].

486    Predictions for translation initiation codons from the *FMR1* sense strand upstream from the

487    repeat identify nine codons that could be used to initiate translation of poly-G, and two for poly-

488    R. The predicted GTG located 11 bases upstream, the ACG located 35 bases upstream, and ACG

489    located 60 bases upstream, have been confirmed experimentally [36]. The antisense upstream

490    region has a total of sixteen codons predicted to initiate translation of poly-P, three for poly-R,

491    and one for poly-A.

492    For the *DM1* sense strand upstream from the repeat, the software predicts three codons that

493    initiate translation of poly-C, and two that initiate translation of poly-A. Interestingly, every CTG

494    within the CTG repeat expansion is predicted to initiate translation of poly-L; however, only the

495    first has a relatively high KSS (0.67). Predictions for the DM1 antisense strand are different from

496    those produced for the experimentally modified DM1 antisense strand (Table 1). Namely, there

497    is no predicted ATG located 17 bases upstream of the repeat expansion, nor a predicted ATT

498    located 23 bases upstream of the repeat expansion, since sequences that border the predicted

499    codons in the modified strand differ from those bordering the same codons in the unmodified

500    version. In the unmodified antisense strand, there are seven codons predicted to initiate poly-A

501    translation, and one to initiate translation of poly-S. Also, there are no predicted translation

502     initiation codons in the reading frame of poly-Q which suggests that this polypeptide might be

503     initiated from the repeat expansion, possibly by repeat length-dependent folding.

504     With respect to the *HDL2* sense strand upstream from the repeat, the software predicts seven

505     codons to initiate translation of poly-L, one to initiate translation of poly-C, and two to initiate

506     translation of poly-A. Furthermore, the software suggests that every CTG of the CTG repeat

507     expansion, aside from the first one in the sense strand, can initiate translation of poly-L. In the

508     antisense strand, there are seventeen codons predicted to initiate translation of poly-Q, three for

509     poly-S, and two for poly-A. The predicted ATC located 6 bases upstream of the repeat expansion

510     in the antisense strand has been confirmed [35].

511     Predictions for translation initiation codons from the *HTT* sense strand upstream from the repeat

512     identify seventeen codons that initiate translation of poly-Q, and four for poly-A. From the

513     antisense upstream region, sixteen codons are predicted to initiate translation of poly-L, and nine

514     for poly-A. The software also suggests that every CTG of the CTG repeat expansion, aside from

515     the first one in the antisense strand can initiate translation of poly-L.

516     Predictions for the *DM2* sense strand upstream from the repeat identify five codons used for

517     translation initiation of poly-PACL, two for poly-CLPA, and three for poly-LPAC. Moreover,

518     the software predicts the first two CTGs of the CCTG repeat expansion to initiate translation of

519     poly-LPAC. In the antisense strand, there are three codons predicted to initiate poly-RQAG

520     translation, five to initiate translation of poly-GRQA, and one to initiate translation of poly-

521     QAGR.

522

523

# Discussion

524

525 As shown here, RFCs were able to successfully predict most translation initiation codons

526 associated with neurologic repeat expansion diseases that were experimentally identified. The

527 same models also predicted other codons to initiate translation of repeat expansions for

528 neurologic diseases, that have not been identified. Of note, this software predicted translation

529 initiation sites with more than 18% accuracy than the TITER neural network.

530 Regardless of the quality of a model, its predictions should not be interpreted as evidence.

531 Instead, predictions should be recognized as likely possibilities that warrant further investigation.

532 The significance of the algorithm's identification of translation initiation codons, however,

533 should not be understated. For example, these data may be important to use to guide treatment of

534 these repeat diseases.

535 Although the machine learning models show promise in understanding of the pathogenesis of

536 repeat expansion neurologic disorders, their use may be extended to other applications as well.

537 For example, they may be used to predict the translation initiation that are not involved in repeat

538 expansion disorders. One benefit of this implementation includes the ability to speculate protein

539 products from a nucleotide sequence, quickly and easily and without laboratory procedures. In

540 order to accelerate the use of the RFCs, a version of the machine learning software that can

541 predict translation initiation codons in any provided sequence is available (at

542 www.tispredictor.com/tis).

543

## Enhancing Performance

545 Like other machine learning models, RFC performance is determined by the amount of available

546 training data. Because of this constraint, collecting more examples to train the machine learning

547 models could prove especially useful. In the case of the near-cognate RFC, obtaining sufficient

548 data to account for all near-cognate types could lessen uncertainty in predictions involving these

549 codons. Training the two RFCs discussed here with more of the codon types that have been used

550 would be beneficial since feeding a model with more data will verify existing trends, and

551 introduce variations that the algorithm can recognize and link to a particular classification,

552 thereby improving accuracy.

553

554

# Materials and Methods

## Data acquisition

557 Examples of translation initiation were mostly obtained from ribosome profiling, mass

558 spectroscopy, and CRISPR-based techniques across different human cell types and under

559 different conditions [38]. These data include sequences of 12,094 examples of translation

560 initiated from ATG, as well as 2,180 examples of translation initiated from near-cognate codons.

561 Translation initiation sites were also captured by quantitative translation initiation sequencing of

562 genes in cultured human kidney cells [39]. Their annotated sequences were collected from the

563    Ensembl gene annotation system (version 84) [40]. These methods procured 509 and 203 more

564    examples of ATG and near-cognate initiation codons, respectively. In all, we collected 12,603

565    instances of translation initiation from ATG, and 2,413 instances of translation initiation from

566    near-cognate codons to use in this study.

567    To obtain examples in which translation does not initiate from ATG (negatives), we used the

568    same transcripts from which positives were derived and recorded nucleotides that flanked ATG

569    codons. Then, we eliminated all instances in which flanking sequences matched any of the

570    12,603 sequences bordering the known ATG translation initiation sites, leaving 34,097

571    negatives. We repeated the same procedure to identify negatives for near-cognate codons that do

572    not initiate translation. We found examples of CTG, GTG, and TTG codons in which flanking

573    sequences did not match any of that of the known near-cognate initiation codons, leaving

574    141,071 negatives.

575

## Random Sampling

577    All random sampling was conducted without replacement. This method is preferred for KSS

578    evaluations of ATG and near-cognate codons, as the precision of population estimates is higher

579    than that produced by sampling with replacement [41]. Furthermore, sampling without

580    replacement to generate training datasets introduces greater variation for model training.

581

## Random forest classifiers

583    Using the open-source package, imbalanced-learn, in Python, we created the RFC models [33].

584    The ATG RFC was trained on an imbalanced set of 12,432 ATG codons known to initiate

585 translation (positives), and 3,261 ATG codons that are believed not to initiate translation

586 (negatives). The set of 3,261 negatives consisted of 1,716 sequences that were not missing

587 nucleotides, and 1,545 (ten percent fewer) randomly sampled negatives of the remaining 31,697

588 that were missing nucleotides. To clarify, missing nucleotides are registered in the case that a

589 recorded codon is located exceedingly close to the 5' or 3' terminus of an mRNA construct. In

590 such a circumstance, there may not be a full ten bases both preceding and following the codon.

591 The sampling technique was performed to slightly offset the proportion of negatives with and

592 without missing bases in the opposite direction. In this way, more negatives without missing

593 bases would be used for model training. Using the original imbalanced set of negatives, with the

594 majority missing bases, would cause the model to inaccurately assess the effect of missing

595 nucleotides on a codon's ability to initiate translation. Furthermore, using a slightly larger

596 proportion of negatives that had a complete sequence profile resulted in improved accuracy for

597 distinguishing codons that were not missing nucleotides. This is useful, as sequences are less

598 often encountered with missing nucleotides in real-world applications.

599 To account for the imbalance of positives and negatives, the RFC had decision trees generated

600 from 3,576 negatives, and the same number of randomly sampled positives. One thousand such

601 trees were used, since this number is generally recommended as a starting point for the

602 generation of an RFC [42]. Of the total number of features, $n$, a total of $\sqrt{n}$ features were used to

603 classify the data in order to optimize predictive power. Training with too many or too few

604 features could have prevented the model from recognizing the best indicators for classification

605 [42]. Each decision tree also had the requirement of grouping at least two codon instances to a

606 certain classification. This constraint reduced the risk of overfitting, yet still allowed tree

607 capacity to differentiate between subtly differing codons. Thus, the trees could better identify

608    precise feature patterns to associate with a particular classification, and remain reliable in face of

609    new, unencountered data.

610    We evaluated the accuracy of the RFC model with the above configurations. Parameters such as

611    the minimum number of codons to group for classification could then be adjusted to improve

612    predictive power, as necessary. However, parameters were best left unchanged for optimal

613    predictions. To create a separate classifier for near-cognate codons, we repeated the same

614    procedures to create an RFC for near-cognate codons as we had carried out for ATG codons, this

615    time using data available for all near-cognate codons.

616

## 617    Accessibility and implementation

618    The software is publicly accessible as an interactive website at www.tispredictor.com.

619    [Availability Statement for Open Access Models and Data]

620

621

622

623

624

625

626

# References

1.    La Spada AR, Taylor JP. Repeat expansion disease: Progress and puzzles in disease pathogenesis. Nat Rev Genet. 2010;11(4):247-58.

2.    Davis M, Stroud C, editors. Neurodegeneration: Exploring Commonalities Across Diseases: Workshop Summary. Forum on Neuroscience and Nervous System Disorders; 2013; Washington (DC): National Academies Press (US).

3.    Rudich PD, Watkins S, Lamitina T. PolyQ-independent toxicity associated with novel translational products from CAG repeat expansions. PLoS One. 2020;15(4).

4.    Boivin M, Pfister V, Gaucherot A, Ruffenach F, Luc N, Sellier C, et al. Reduced autophagy upon C9ORF72 loss synergizes with dipeptide repeat protein toxicity in G4C2 repeat expansion disorders. EMBO J. 2020;39(4):e100574.

5.    Boivin M, Deng J, Pfister V, Grandgirard E, Oulad-Abdelghani M, Mortlet B, et al. Translation of GGC repeat expansions into a toxic polyglycine protein in NIID defines a novel class of human genetic disorders: The polyG diseases. Neuron. 2021;109(11):1825-35.

6.    Lee S-J, Hee-Sun L, Masliah E, Lee H-J. Protein aggregate spreading in neurodegenerative diseases: Problems and perspectives. Neurosci Res. 2011;70(4):339-48.

7.    Monaco A, Fraldi A. Protein Aggregation and Dysfunction of Autophagy-Lysosomal Pathway: A Vicious Cycle in Lysosomal Storage Diseases. Frontiers in Molecular Neuroscience. 2020.

8.    Ross CA, Poirier MA. Protein aggregation and neurodegenerative disease. Nat Med. 2004;10 Suppl:S10-7. Epub 2004/07/24. doi: 10.1038/nm1066. PubMed PMID: 15272267.

648    9.    Chung CG, Lee H, Lee SB. Mechanisms of protein toxicity in neurodegenerative

649    diseases. Cellular and Molecular Life Sciences. 2018;75(17):3159-80. doi: 10.1007/s00018-018-

650    2854-4.

651    10.    Krans A, Skariah G, Zhang Y, Bayly B, Todd PK. Neuropathology of RAN translation

652    proteins in fragile X-associated tremor/ataxia syndrome. Acta Neuropathologica

653    Communications. 2019;7(1).

654    11.    Zhang S, Hu H, Jiang T, Zhang L, Zeng J. TITER : predicting translation initiation sites

655    by deep learning. Bioinformatics. 2017:33:i234–i42.

656    12.    Reuter K, Biehl A, Koch L, Helms V. PreTIS: A Tool to Predict Non-canonical 5' UTR

657    Translational Initiation Sites in Human and Mouse. PLOS Computational Biology.

658    2016;12(10):e1005170. doi: 10.1371/journal.pcbi.1005170.

659    13.    Hernández G, Osnaya VG, Pérez-Martínez X. Conservation and Variability of the AUG

660    Initiation Codon Context in Eukaryotes. Trends in Biochemical Sciences. 2019;44(12):1009-21.

661    14.    Meijer HA, Thoma AAM. Control of eukaryotic protein synthesis by upstream open

662    reading frames in 5′-untranslated region of an mRNA. Biochem. 2002;367:1-11.

663    15.    Pisarev AV, Kolupaeva VG, Pisareva VP, Merrick WC, Hellen CUT, Pestova TV.

664    Specific functional interactions of nucleotides at key −3 and +4 positions flanking the initiation

665    codon with components of the mammalian 48S translation initiation complex. Genes Dev.

666    2006;20:624-36.

667    16.    Lütcke HA, Chow KC, Mickel FS, Moss KA, Kern HF, Scheele GA. Selection of AUG

668    initiation codons differs in plants and animals. Embo J. 1987;6(1):43-8.

669    17.    Kozak M. At least six nucleotides preceding the AUG initiator codon enhance translation

670    in mammalian cells. J Mol Bio. 1987;196(4):947-50.

671    18.    Kozak M. Recognition of AUG and alternative initiator codons is augmented by G in

672    position +4 but is not generally affected by the nucleotides in positions +5 and +6. Embo J.

673    1997;16(9):2482-92.

674    19.    Kozak M. Point mutations define a sequence flanking the AUG initiator codon that

675    modulates translation by eukaryotic ribosomes. Cell. 1986;44(2):283-92.

676    20.    Kozak M. Context effects and inefficient initiation at non-AUG codons in eucaryotic

677    cell-free translation systems. Mol Cell Biol. 1989;9(11):5073-80.

678    21.    Wei J, Zhang Y, Ivanov IP, Sachs MS. The stringency of start codon selection in the

679    filamentous fungus Neurospora crassa. J Biol Chem. 2013;288(13):9549-62. Epub 2013/02/08.

680    doi: 10.1074/jbc.M112.447177. PubMed PMID: 23396971.

681    22.    Kearse MG, Wilusz JE. Non-AUG translation: a new start for protein synthesis in

682    eukaryotes. Genes Dev. 2017;31(17):1717-31. Epub 2017/10/07. doi: 10.1101/gad.305250.117.

683    PubMed PMID: 28982758; PubMed Central PMCID: PMCPMC5666671.

684    23.    Libbrecht MW, Noble WS. Machine learning in genetics and genomics. Nat Rev Genet.

685    2015;16(6):321-32.

686    24.    Schneider TD, Stephens RM. Sequence logos: a new way to display consensus

687    sequences. Nucleic Acids Research. 1990;18(20):6097-100.

688    25.    Wegrzyn J. L. DTM, Valafar F., Hook V. Bioinformatic analyses of mammalian 5'-UTR

689    sequence properties of mRNAs predicts alternative translation initiation sites. BMC

690    Bioinformatics. 2008.

691    26.    Schwab SR, Shugart JA, Horng T, Malarkannan S, Shastri N. Unanticipated Antigens:

692    Translation Initiation at CUG with Leucine. PLoS Biology. 2004;2(11):e366.

693    27.    Davis J, Goadric M. The relationship between Precision-Recall and ROC curves. ICML.

694    2006:233-40.

695    28.    Saito T, Rehmsmeier M. The Precision-Recall Plot Is More Informative than the ROC

696    Plot When Evaluating Binary Classifiers on Imbalanced Datasets. PLOS ONE. 2015.

697    29.    Libbrecht M. W. NWS. Machine learning in genetics and genomics. Nat Rev Genet.

698    2015;16(6):321-32.

699    30.    Goldstein BA, Polley EC, Briggs FBS. Random Forests for Genetic Association Studies.

700    Stat Appl Genet Mol Biol. 2011;10(1):32.

701    31.    Liu Y, Wang Y, Zhang J. New Machine Learning Algorithm: Random Forest. ICICA.

702    2012;7473:246-52.

703    32.    Friedman JH. On Bias, Variance, 0/1—Loss, and the Curse-of-Dimensionality. Data

704    Mining and Knowledge Discovery volume. 1997;1(1):55–77.

705    33.    Lemaître G, Nogueira F, Aridas CK. Imbalanced-learn: A Python Toolbox to Tackle the

706    Curse of Imbalanced Datasets in Machine Learning. JMLR. 2017;18(17):1-5.

707    34.    Fawcett T. ROC Graphs: Notes and Practical Considerations for Researchers. Pattern

708    Recognition Letters. 2004;31(8):1-38.

709    35.    Zu T, Gibbens B, Doty NS, Gomes-Pereira M, Huguet A, Stone MD, et al. Non-ATG-

710    initiated translation directed by microsatellite expansions. Proc Natl Acad Sci U S A.

711    2010;108(1):260-5.

712    36.    Rodriguez CM, Wright SE, Kearse MG, Haenfler JM, Flores BN, Liu Y, et al. A native

713    function for RAN translation and CGG repeats in regulating fragile X protein synthesis. Nature

714    Neuroscience. 2020;23:386–97.

715    37.    Kearse MG, Green KM, Krans A, Rodriguez CM, Linsalata AE, Goldstrohm AC, et al.

716    CGG Repeat associated non-AUG translation utilizes a cap-dependent, scanning mechanism of

717    initiation to produce toxic proteins. Mol Cell. 2016;62(2):314-22.

718    38.    Chen J, Brunner A-D, Cogan JZ, Nunez JK, Fields AP, Adamson B, et al. Pervasive

719    functional translation of noncanonical human open reading frames. Science. 2020;367:1140–6.

720    39.    Gao X, Wan J, Liu B, Ma M, Shen B, Qian S-B. Quantitative profiling of initiating

721    ribosomes in vivo. Nature Methods. 2015;12(2):147-53.

722    40.    Aken BL, Ayling S, Barrell D, Clarke L, Curwen V, Fairley S, et al. The Ensembl gene

723    annotation system. Database (Oxford). 2016.

724    41.    Seth GR, Rao JNK. On the Comparison between Simple Random Sampling with and

725    without Replacement. Sankhyā: The Indian Journal of Statistics, Series A (1961-2002).

726    1964;26(1):85-6.

727    42.    Johnson K, Kuhn M. Applied Predictive Modeling. New York: Springer; 2013.

Figure 1

Kozak Similarity Scores of Known ATG Translation Initiation Codons Against Baseline

Legend:
- 12,604 Known ATG Translation Initiation Codons
- 100,000 Codons With Randomized Flanking Sequences
- 34,097 ATG Codons Believed Not To Inititate Translation

X-axis: Kozak Similarity Score of Codon

Y-axis: Proportion

Figure 2

**Kozak Similarity Scores of Known Near-Cognate Translation Initiation Codons Against Baseline**

Legend:
- 1,514 Known CTG Translation Initiation Codons
- 543 Known GTG Translation Initiation Codons
- 356 Known TTG Translation Initiation Codons
- 100,000 Codons With Randomized Flanking Sequences
- 141,071 CTG, GTG, and TTG Codons Believed Not to Inititate Translation

Figure 3

Figure 4

Figure 5

## Proportion of Codons that Initiate Translation With Kozak Similarity Scores Above Certain Values

| Kozak Similarity Score | Proportion of ATGs that Initiate Translation above Kozak Similarity Score | Proportion of Near-Cognate Codons that Initiate Translation Above Kozak Similarity Score |
|---|---|---|
| **0.00** | 0.5000 | 0.5000 |
| **0.05** | 0.5000 | 0.5000 |
| **0.10** | 0.5000 | 0.5000 |
| **0.15** | 0.5001 | 0.5000 |
| **0.20** | 0.5033 | 0.5035 |
| **0.25** | 0.5098 | 0.5110 |
| **0.30** | 0.5206 | 0.5224 |
| **0.35** | 0.5375 | 0.5412 |
| **0.40** | 0.5618 | 0.5683 |
| **0.45** | 0.5933 | 0.6015 |
| **0.50** | 0.6302 | 0.6406 |
| **0.55** | 0.6789 | 0.6893 |
| **0.60** | 0.7364 | 0.7469 |
| **0.65** | 0.7938 | 0.8011 |
| **0.70** | 0.8496 | 0.8541 |
| **0.75** | 0.8948 | 0.9067 |
| **0.80** | 0.9393 | 0.9458 |
| **0.85** | 0.9682 | 0.9753 |
| **0.90** | 0.9908 | 0.9779 |
| **0.95** | 1.0000 | No codons above score |

Figure 6

Figure 7

Receiver Operating Characteristic Curves of Codon Classifiers (Derived from Same Data)

Figure 8

TACCTTGTAGAAAGCGCCATTGGAGCCCCGCACTTCCACCACCAGCTCCTCCATCTTCTCTTCAGC
CCTGCTAGCGCCGGGAGCCCGCCCCCGAGAGGTGGGCTGCGGGCGCTCGAGGCCC

Figure 9

**C9orf72 (Sense):**
Reading Frame: **Poly-GA** = blue, **Poly-GR** = purple

**CTG**GAACTCAGGAGTCGCGCGCTA**GG**GGCC(**GGGGCC**)ₙ*

24 ↑          ↑
           1

**C9orf72 (Antisense)**
**Poly-PR** reading frame = **blue**, **Poly-PG** reading frame = **purple**

**AAG**CCGCGCGCCGCCCACCCTCCGGCCTTCCCCCAGGCG**AGG**CCTCTCAGTACCCGAGGCTCC
492 ↑                                    ↑
                                       453

CTTTTCTCGAGCCCGCAGCGGCAGCGCTCCCAGCGGGTCCCCGGG**AAG**GAGACAGCTCGGGTA
                                              ↑
                                            384

**CTG**AGGGCGGGAAAGC**AAG**GAAGAGGCCAGA**TC**CCCATCCCTTGTCCCTGCGCCGCCGCCGCC
↑             ↑              ↑
366          350           336

GCCGCCGCCGCCGGG**AAG**CCCGGGGCCCGG**ATG**CAGGCA**ATT**CCACCAGTCGCTAGAGGCGAA
              ↑              ↑        ↑
             288           273      264

AGCCCGACACCCAGCTTCGGTCAGAGAA**ATG**AGA**GG**GAAAGTAAAA**ATG**CGTCGAGCT**CTG**A**GG**
                           ↑       ↑          ↑              ↑   ↑
                          212     207        194            182 179

AGAGCCCCGCTTCTACCCGCGCCTCTTCCCGGCAGCCGAACCCCAAACAGCCACCCGCCAGG**A**
                                                               ↑
                                                              113

**TG**CCGCCTCCTCACTCACCCACTCGCCACCGCCTGCGCCTCCGCCGCCGCGGGCGCAGGCACC

GCAACCGCAGCCCCGCCCCGGGCCCGCCCCGGGCCCGCCCCGACC**ACG**(**CCCCGG**)ₙ
                                            ↑

**FMR1 (Sense)**
Reading Frame: **Poly-G** = blue, **Poly-R** = purple

**CTG**AGTGCACCT**CTG**CAGAAATGGGCGTTCTGGCCCTCGCGAGGCAGTGCGACCTGTCACCGCC
↑           ↑
419        407

CTTCAGCCTTCCCGCCCTCCACC**AAG**CCCGCGCACGCCCGGCCCGCGCGT**CTG**TCTTTCGACCC
                      ↑                            ↑
                     332                          305

GGCACCCCGGCCGGTTCCCAGCAGCGCGCATGCGCGCGCTCCCAGGCCACTTGAAGAGAGAGG

GCGGGGCCGAGGGGCTGAGCCCGCGGGGGGGAGGGAACAGCGTTGATCACGTGA**CG**TGGTTTCA
                                                      ↑
                                                     176

GTGTTTACACCCGCAGCGGGCCGGGGGTTCGGCCTCAGTC**AGG**CGCTCAGCTCCGTTTCGGTTT
                                        ↑
                                       125

CACTTCCGGTGGAGGGCCGCCTCTGAGCGGGCGGCGGGCCG**ACG**GCGAGCGCGGGCGGCGGC
                                         ↑
                                        60

G**GTG**A**CG**GAGGCGCCG**CTG**CC**AGG**GGGC**GTG**CGGCAGCG(**CGG**)ₙ
 ↑   ↑            ↑     ↑      ↑
38  35           23    18     11

**FMR1 (Antisense)**
Reading Frame: **Poly-P** = blue, **Poly-R** = purple, **Poly-A** = green

**CTG**CCGCCGGCCCTCGCCCATCCCCAGCTCACCCCGGCGGGGCTCGGCGCCGAAAGAGAACCC
↑
414

TCTCCTCGCTGGTCTCTCATTTCGATAGGCGCTA**GG**GCCTCTCGGAGTCGGGAGAGGGGCTTCCAA
                                 ↑
                                318

C**AGG**CCCCAAGTCCAGTCCTTCCCTCCCAACAAC**ATC**CCGCCGAGC**GTG**CCCTGGCACCCAGGC
↑                                 ↑            ↑
285                              252          240                  168 ↓

GCGGTGCTCGGG**AAGAGG**GCCCCGGGC**CTG**CCTCCCGCCGACACC**AAGAAG**AAA**AGG**GAGGGA
          ↑     ↑           ↑                ↑    ↑      ↑
         210   207         195              177  174    105

**AGGAAG**GGCGAAGATGGGGCCTGCCCTAGAGCC**AAG**TACCTTGTAGAAAGCGCC**ATTG**GAGCCC
↑    ↑                           ↑                        ↑
159 156                         126                      104

CGCACTTCCACCACCAGCTCCTCCATCTTCTCTTCAGCC**CTG**CTAGCGCCGGGGAGCCCGCCCCC
                                       ↑
   28 ↓                                56

GAGAG**AG**G**GTG**GG**CTG**CGGGCGCTCGAGGCCCAG(**CCG**)ₙ
    ↑    ↑
   26   21

# Figure 10

## DM1 (Sense)
Reading Frame: Poly-C = blue, Poly-A = purple, Poly-L = green

**CTG**CCAGTTCACAACCGCTCCGAGCGTGGGTCTCCGCCCAGCTCCAGTCCTGTG**ATC**CGGGCCC
↑ 143          ↑ 89  ↓ 23

GCCCCCTAGCGGCCGGGG**AGG**GAGGGGCCGGGTCCGCGGCCGGCGAACGGGGGCTCG**A**A**AGG**G
↑ 61          ↑ 22

TCCTTGTAGCCGGGAATG**(CTG)n***

## DM1 (Antisense)
Reading Frame: **Poly-A = blue**, Poly-S = purple

**TTG**GCAAAAGCAAATTTCCCGAGTAAGCAGGCAGAG**ATC**GCGCCAGACGCTCCCCAGAGC**AGG**
↑ 229          ↑ 193          ↑ 169

CGTCATGCACAAGAAAGCTTTGCACTTTGCGAACCAACGATAGGTGGGGGTGCGTGGAGGATGG

AAC**A**AGG**CACCCCCCGCTTCGCTGCCTTCCCA**CCC**CTC**CAGTTTGCCCATCCACGTCAGGGCCTC
↑ 98 ↑ 94          ↑ 67

AGC**CTG**GCCGAAAGA**AAG**AAATGGTCTGTGATCCCCC**(CAG)n**
↑ 34          ↑ 22

## HDL2 (Sense)
Reading Frame: **Poly-L = blue**, Poly-C = purple, Poly-A = green

**CTGCTG**GAGGGAGGAGGG**AGG**CCCATCTGCTCA**GTG**AGAGCCCAGGAATCTCGTCTTTCAGTGG
↑ 165 ↑ 162          ↑ 147          ↓ 68 ↑ 132

CTGCATCGTTTTCACCATTAGTTGAGGGA**ATC**GAT**CTG**TGCCTTCATTCTAAG**ATG**CCACCGCATT
↑ 72          ↑ 66          ↑ 48

CGGGGCAGAGCCGGGGCCGGAAGCC**AGG**GAG**CTG**C**(CTG)n**†
↑ 10          ↑ 4

## HDL2 (Antisense)
Reading Frame: **Poly-Q = blue**, Poly-S = purple, Poly-A = green

**CTG**CCACTCC**CTG**GAT**GTG**GCC**GTG**GAGCAG**GTG**TTCCGAGCTTTCTTCATGAGACTCGGAGGCA
↑ 393          ↑ 384          ↑ 378          ↑ 372          ↑ 363

GCGAGCGGGCGGGG**AGG**ACTAGAGCACACGTACGCTCAGGATTTTTATTCTCGGCAGAAGCACA
↑ 315

AGTC**ATT**GCTGCGTCTCACTCCCTTCCCAGGCATCCATGCCTGAAACAGGCGAGCCGCCCG**AGG**
↑ 261          ↑ 204

GAGGCAGAGGACCCA**GTG**CCTCTGAGG**ATA**AGCTCCCAC**ATG**GGA**GTG**GACTGGTCGGGGAAG
↑ 186          ↑ 174          ↑ 162          ↑ 156

CCAGTTTGTTTGGGGT**CTG**GCTTCCAGGAA**AGG**ACCGGCTCC**ATG**CGTGTC**AGG**TGCA**CTG**AGG
↑ 122          ↑ 108          ↑ 96          ↑ 87          ↑ 80

A**GTG**GAT**ATC**GGAGAGTCCAGGCAGGCGGGCAG**CTG**CAGAGCCGGCC**AAG**GTTCCCTGCACAG
↑ 73 ↑ 67          ↑ 41          ↑ 27

AAACC**ATC**TTA**(CAG)n**
↑ 6

Figure 11

## HTT (Sense)
Reading Frame: **Poly-Q** = blue, **Poly-A** = purple

**CTG**CGCTGTCAGCGGCCTTGCTGTGTGAGGCAGAAC**CTG**CGGGGGCAGGGGCGGGCTGGTTCC
384        348

**CTG**GCCAGCCATTGGCAGAGTCCGCAGGCT**AGG**GCTGTCAATCATGCTGGCCGGCGTGGCCCC
321      291

GCCTCCGCCGGCGCGGCCCCGCCTCCGCCGCGCAGCGT**CTG**GG**ACG**CAAGGCGCC**GTG**GGG
219   214   210   202

G**CTG**CCGGG**ACG**GGTCCAAGATGGACGGCCGCTCAGGTTCTGCTTTTAC**CTG**CGGCCCAGAGCC
195    187            147

CCATTC**ATT**GCCCCG**GTGCTG**AGCGGCGCCGCGAGTCGGCCCGAGGCCTCCGGGGA**CTG**CCGT
126     117   114        76

GCCGGGCGGGAGACCGCC**ATG**GCGACC**CTG**GAAAAG**CTG**ATGAAGGCCTTCGAGTCCCTC**AAG**
51    42   36   33   27      9

TCCTTC**(CAG)ₙ**

## HTT (Antisense)

Reading Frame: **Poly-L** = blue, **Poly-?** = purple

**AGG**CCCCAACAAGGCTCTGCCTCCCCCTCGCGAGAGGACA**AGG**GAAGACCCAAGTGAGGGAGCG
408    339              369

GGG**CTG**AA**GTG**GGGGAAGGCCTCGCCCCAGGAGGGGGCGGGTGTCCCTCATGGGCT**CTG**GGTT
342   337    330           289

GCTGGGTCACT**CTG**TCT**CTG**CGGGGCCGGGGGTTCGTGTCGCCGGCCCGCAGG**CTG**C**AGG**GTT
271   265           229   225

ACCGCC**ATC**CCCGCCGTAGC**CTG**GGACCCGCCGGGACAGGGAGCTGCAGCGGGCCCAAACTCA
213     199

CGGTCGGTGCAGCGGCTCCTCAGCCACAGCCGGGCCGG**GTG**GCGGCGGGGGCGGCGGCGGG
118

GGCGG**CTG**CGG**CTG**AGGCAGCAGCGG**CTG**TG**CCTGCGGCGGCGGG**CTG**AGGAAG**CTG**AGGAGG
90     84       69 67        51       42   39   36

CGGCGGCGGCGGCGGCGGCG**GTG**GCGG**CTG**TTG**(CTG)***
13     6

## DM2 (Sense)
Reading Frame: **Poly-PACL** = blue, **Poly-CLPA** = purple, **Poly-LPAC** = green

**ATT**ACTGCCAGT**GTG**TGTGTGTGTGTGTGTGTGTGTGTGTGTGTGTGTCTGTCTGT**CTG**T**CTG**TCT
87      75                     31     27     23

G**TCTGTCTGTCTGTCTGTCTGCCTGCCTG(CCTG)ₙ†**
19    15    11     7     3     -2     -6

## DM2 (Antisense)
Reading Frame: **Poly-RQAG** = blue, **Poly-GRQA** = purple, **Poly-QAGR** = green

**GTG**GAGGTTGCA**GTG**AGCCGAGAT**CATA**CCA**CTG**CACTCCAGCCTAGGGGACAAA**GTG**AGACAG
113       101       91   88     82            58

AC**AGG**CAGGCAGGCAGGCAGGCAGGCAGGCAGAC**AGG**CAGACAGGCAGC**(CAGG)ₙ**
47                    15     7

# Figure 12