
Generating Football Match Summaries with NMT

Miguel Ayala

Department of Computer Science
Stanford University
maya.la3@stanford.edu

Abstract

1 In sports journalism, football (soccer) match summaries are handwritten, despite
2 their formulaic nature. There seem to be major patterns that govern this style of
3 writing. Given the improvement in text generation algorithms, this study aims to
4 see whether or not football match summaries can be produced with neural network
5 models. An NMT framework is used to generate textual summaries from game
6 event vectors. For this reason and several others, BLEU scores remain low. While
7 semi-coherent text is produced by the neural network, it lacks specificity regarding
8 match in question. What is promising, however, is that the generated text seems
9 to be salient to some degree. In other words, if the game event vector described a
10 goal, then the generated summary was often about a goal-scoring chance. Future
11 work will require specific entity handling to ensure that generated sentences are
12 effective.

13 **1 Introduction**

14 **1.1 The Problem**

15 With so much demand for sports content in this age of digital consumption, there have been numerous
16 attempts to procedurally generate summaries of sporting events. Football (soccer) still relies on
17 journalists to cover the thousands of matches that happen every week. Demand for instant summaries
18 is incredibly high, however.

19 **1.2 The Goal**

20 The goal of this study is to produce a neural network that can produce instant textual summaries of
21 football matches. More concretely, the objective is to take a vector of comprehensive match event
22 data and generate a coherent narrative that mentions only the salient events in a fluid manner.

23 **1.3 Related Work**

24 Different approaches seem to work for other sports. Successful examples include the WordSmith
25 program which, among other things, produces basketball game summaries for the Associated Press
26 [7]. As of 2019, there is no such parallel for football (soccer) matches. Nevertheless, there have
27 been attempts over the past 3 decades. Each approach is fundamentally different. The input data,
28 for instance, is incredibly varied. André et al tried to transform visual images of football matches
29 into text as far back as 1988 [1]. In 1998, abstracted representations were being used for soccer
30 commentary generation [3]. In recent years, however, many studies have focused on using game
31 statistics to create match summaries[4] [5]. These approaches are still largely rule-based, however,
32 and the summaries produced fail to capture the fluidity and flair of human written summaries.

33 One framework that might be able to provide this level of linguistic flexibility is the NMT [9]. NMT
34 models are largely based on the encoder-decoder paradigm [6] [14]. This is in contrast to most
35 traditional, phrase translation models that consist of many individual sub-components that have to
36 be trained and tuned separately. In NMT, the encoding phase consists of a bidirectional Recurrent
37 Neural Network that encodes a vector representation of the source sentence. The decoder then takes
38 this encoder output to produce a target sentence. This relatively simple process now underpins some
39 of the most effective translation services currently available[15] [8]. NMT combines the strengths of
40 word embeddings with memory models to create a system that considers entire sentences and learns
41 deep links between words.

42 Most of the current research still lies around the original task of translation. Different institutions are
43 tweaking the various components of the basic encoder-decoder model. Some are trying to work on
44 problems like the fixed vector length problem [2] or the vanishing/exploding gradient problem. Other
45 optimization comes in the various attention mechanisms that can be used to condense or remember
46 lengthy source vectors [10]. While some applications exist, there is still potential for novel NMT
47 experimentation. This problem may be one such use case.

48 **2 Approach**

49 **2.1 The NMT Approach**

50 The task of producing match summaries from statistics is not explicitly a translation task. It does not
51 map strings to strings. Instead it maps numbers to strings. The NMT is flexible enough to handle this,
52 however. Instead of passing a word-embedding representation of a source sentence to the decoder, we
53 pass a mix of vectorized match statistics and word-embeddings, with appropriate padding to ensure
54 uniform length. To be more specific, our input is a flattened vector consisting of event vectors ¹.

55 **2.2 Architecture**

56 The NMT system's architecture is inspired by one built for Winter 2019 Natural Language Processing
57 with Deep Learning at Stanford (see Figure 1). The system consists of a bidirectional LSTM encoder
58 and a unidirectional LSTM decoder. Words in a source sentence are converted into its embeddings.
59 These are passed through to the encoder in exchange for hidden and cell states. The decoder's
60 LSTM is initialized with hidden and cell states from this process. Multiplicative attention is used in
61 conjunction with a linear layer, a tanh layer and a dropout layer to produce a combined output vector.
62 From this vector we can produce a probability distribution of likely target words. From this, we can
63 generate likely sentences.

¹see Data for more details

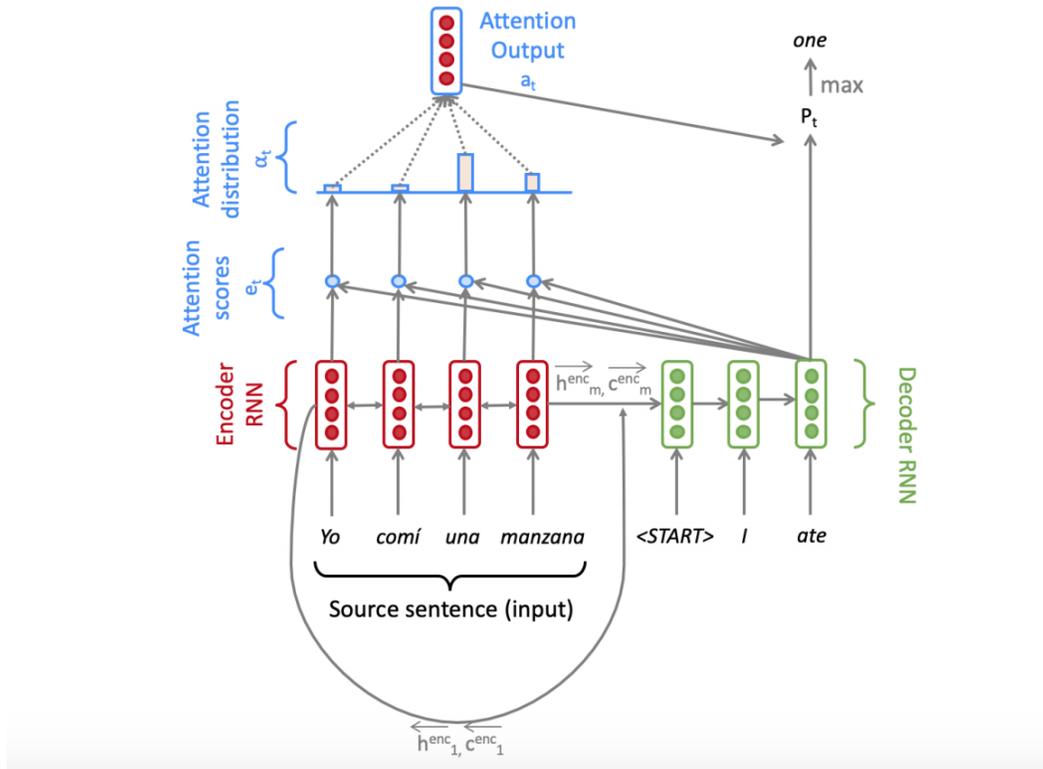


Figure 1: Architectural Diagram

64 3 Data

65 3.1 Data Collection

66 The dataset is bipartite. It consists of source vectors and targets. Each was collected from a different
 67 source. A variety of preprocessing tools were built to get the data into usable format and then to form
 68 associations between source sentences and target sentences.

69 3.1.1 Source Vectors

70 The source vectors come from a curated dataset [12] that details major football events in matches
 71 from the 5 biggest European leagues: The English Premier League, La Liga, Bundesliga, Serie A and
 72 Ligue 1. In contrast to many football match datasets, each datum from the Football Events dataset
 73 consists of far more than just the typical goals, red cards, yellow cards, etc. The data for each match
 74 is broken down into key events that are quantified by a series of variables such as *shot_outcome* and
 75 *time*. Here is the full list of categories:

76
 77 ["time", "event_type", "event_type2", "side", "event_team", "opponent", "player", "player2",
 78 "player_in", "player_out", "shot_place", "shot_outcome", "is_goal", "location", "bodypart",
 79 "assist_method", "situation", "fast_break"]

80
 81 Some of these variables are binary or numeric. See Appendix 1 for a legend.

82 3.1.2 Target Vectors

83 The target vectors consist of English language full text match summaries from Sky Sports[13]. This
 84 site was chosen because the style of summary appears to be generally consistent from year to year

85 and from match to match. Additionally, they have match summaries from the Top 5 leagues covered
86 by the Football Events dataset.

87 **3.2 Data Preprocessing**

88 From the Football Events dataset, over 10,113 matches were found. A BS4 scraping tool was
89 designed to pull summaries from SkySports league results pages. These summaries were matched
90 with corresponding leagues and the corresponding seasons. Football Events information about team
91 names, match scores and match dates were used to find the appropriate match summaries. 2,011
92 relevant match summaries were downloaded by the scraper. This gave us a total of 30,712 match
93 summary sentences.

94 The next task was to connect individual game event objects with sentences in these summaries. There
95 were 941,010 game events in total. I developed a regex tool that allowed me to assess the relevance of
96 a game event to a specific summary sentence. This was based on information about players involved
97 and the event described. Through this approach, 12,673 training examples were found. For reference
98 a training example looked like this:
99

100 **Source Vector:**

101

102 [*7 8 na 1 lyon st etienne yoann gourcuff na na na na na 0 2 na 0 na 0*]

103

104 **Target Vector:**

105

106 [*a second half yoann gourcuff strike cancelled out kurt zoumas opener for les verts as lyon missed*
107 *out on a third successive league win but stayed on course for champions league qualification*]
108

109 Each sentence was standardized with case and punctuation. Also, I had to create a dictionary of team
110 names because SkySports and FootballEvents referred to the same team with different names. This
111 allowed for easy cross-referencing between source and target sentences.

112 I then split the set of all source-target pairs into a training set, a development set and a test set of
113 8871, 2535 and 1267 examples respectively.

114 **4 Experiments**

115 **4.1 Evaluation Method**

116 It is unfair to compare the results of this experiment to some of the aforementioned rules-based
117 approaches to match summary generation. In those cases, coherence within a sentence is likely
118 not to be an issue because each sentence is built from a set of prefabricated linguistically correct
119 templates. With the NMT approach in this study, however, it is highly likely that the test output will
120 feature incoherent, nonsensical sentences. Thus, the evaluation method for this generator needs to be
121 different.

122 **4.1.1 BLEU Score**

123 BLEU score is computed [11]. BLEU is useful because it is objective/deterministic and allows for a
124 quick analysis of many test examples. Additionally, BLEU is good because it is a quantified metric
125 meaning that comparison between various translations is easy.

126 To calculate BLEU a few steps are required. First, one must calculate the n -gram precision for
127 each translation. This is given by:
128

$$129 p_n = \frac{\sum_{ngram} \min(\max(Count_{r_i}(ngram), Count_c(ngram)))}{\sum_{ngram} Count_c(ngram)}$$

130

131 Once we have this, we calculate brevity penalty (BP). BP is equal to 1 if $c \geq r^*$ where r^* is the
132 length of the closest reference translation. Otherwise, BP is

133

$$134 \exp\left(\frac{c}{r^*} - c\right)$$

135

136 Then we calculate BLEU with BP and p_n :

137

$$138 \text{BLEU} = \text{BP} \times \exp\left(\sum_{n=1}^4 \lambda_n \log(p_n)\right)$$

139 4.1.2 Subjective Score

140 BLEU also has its drawbacks. Closeness to a reference translation, is not the only goal of a translation
141 system. A relevant and coherent sentence can be a good translation even if it does not closely resemble
142 the reference translation. It might be different just because the reference set is not comprehensive
143 enough. Recognizing this, I devised a small test for how effective each of my outputs was. I asked
144 10 participants to grade a random sample of 20 generated sentences for coherence, specificity (what
145 statistics or names are mentioned?) and saliency (with regards to the context - i.e. do they talk about
146 a goal scoring chance). Each of these was on a scale of 0 to 10. Each person has a different scale
147 but with the random samples and the different viewpoints, the scores might provide insight into the
148 effectiveness of various tests.

149 4.2 Experimental Details

150 I performed 3 main tests. I trained each one on batches of 40. The dropout I used was 0.3. Validation
151 was performed every 200 iterations for maximum of 30 epochs and 5 patience rounds.

152 4.2.1 Naive

153 The first test was a naive approach. I flattened all the event objects into a single vector. This large
154 vector became the source sentence for the NMT system. The target sentence that matched this source
155 sentence was the entire match summary that corresponded to the game in question.

156 4.2.2 Alpha

157 The second experiment consisted of a source sentence that only had the alphabetical components
158 of the event object. For instance, the indicator variables were dropped from the object, leav-
159 ing only player name, team name and a few others. The target sentence was a regex-matched sentence.
160

161 4.2.3 Alpha-Numeric

162 The third test was the approach outlined in **Data**. Each training example consisted of a condensed
163 Football Events event object as the source sentence (alpha-numeric) with a regex-matched sentence
164 as the corresponding target sentence.

165 4.3 Results

166 4.3.1 Evaluation Metrics

Tests			
Metric	Naive	Alpha	Alpha Numeric
167 BLEU	0.163	0.321	0.678
Coherence	3.4	3.8	7.8
Specificity	0.4	0.9	0.5
Saliency	3.3	3.2	4.9

168 Below are some outputs from the Alpha-Numeric test

169 4.3.2 Sample Output 1

170 **Source:**

171 [4 1 na 1 elche getafe hinestroza przemyslaw na na na na 1 4 2 0 1 0]

172 **Reference:**

173 [at the other end on loan midfielder hinestroza cut in from the left but his strike was straight at elche
174 custodian przemyslaw tyton]

175 **Output:**

176 [the spaniard picked up the ball from the edge of the half hour mark but his shot was blocked by
177 lorient]

178 4.3.3 Sample Output 2

179 **Source:**

180 [25 1 12 1 elche getafe rodriguez domingo cisma na na 8 2 0 15 1 1 1 0]

181 **Reference:**

182 [guaita was kept busy at the start of the second half as he had to tip over domingo cismas header as
183 elche picked up where they left off]

184 **Output:**

185 [but came <unk> through the first half with alexis sanchez hitting the bar after <unk> minutes with a
186 lobbed effort from 18 yards]

187 4.3.4 Sample Output 3

188 **Source:**

189 [52 7 na 2 genoa juventus na na eugenio lamanna mattia perin na na 0 na na 0 na 0]

190 **Reference:**

191 [they took the lead in the 25th minute when tevez showed his fine goal past roman weidenfeller on
192 wednesday was no fluke with a similarly stunning strike which hit the underside of the crossbar on its
193 way past mattia perin]

194 **Output:**

195 [despite the other end of <unk> <unk> in the <unk> minute with a brilliant free kick while luis
196 suarez scored twice in the second half to move onto <unk> pass]

197 5 Analysis

198 Because we cannot compare to the rules-based summary generators, we have no real baselines to
199 compare to. The 3 tests do provide interesting points of comparison, however. It is incredibly
200 clear which training paradigm works the best. The Alpha-Numeric approach bested the other two
201 approaches in 3 out of the 4 metrics. Even if the NMT system is not specifically designed to
202 handle numeric inputs, the extra information does seem to improve the performance of the summary
203 generation according to both objective and subjective measures. The only measure where it fails
204 is specificity. All of the other metrics failed at this too, however. That is a failure on the model
205 itself. Despite my best efforts, the model is not equipped to handle non-vocab entities. A rules based
206 approach would probably work in the future.

207 In the other metrics, the approach scores relatively highly. The good coherence and saliency scores
208 are reflected in the sample outputs. In sample output 1, we have the minor non-logical reference to
209 the 'half-hour mark' as a physical location. Nevertheless, it details a shot that is saved by a goalie.
210 This matches both the source vector and the reference. The specificity is non-existent as every name
211 is butchered and replaced with something else.

212 Sample output 2 shows similar strengths and weaknesses. The model clearly understands the context
213 here. The source and reference refer to a missed shot. The output, too, talks about a missed shot. The
214 same mistakes as in output 1 are seen here. All the team names and player names are displaced or
215 missing. We all see the presence of *UNK* tokens.

216 Sample output 3 indicates some larger structural issues. The output is moderately coherent, yet it
217 does not capture the context of the situation and it does not mention any of the relevant specific
218 entities. More concerning, however, is the link between reference and source in this case. These 2
219 sentences do not match up. If inconsistencies like this exist in the training set, then it is very hard to
220 create an effective neural network-based summary generator. Additionally, the reference sentence
221 details events that are outside the scope of the event data. This shows a huge problem in the approach.
222 So much good summary writing takes into account the larger macro context of individual moments.
223 This is something that is not captured by the event object and thus cannot be replicated by the NMT.

224 **6 Conclusion**

225 It is clear that sports journalists do not have to fear. Automation is still a while away. While the
226 results are not perfect, it is a promising start. It seems like a novel way to use NMT frameworks
227 for tasks removed from the traditional sphere of machine translation. If anything, it illustrates how
228 versatile and robust these systems are.

229 There are nonetheless a litany of issues that need to be sorted out if any progress is to be made.

230 The model is currently not equipped to handle numerical information that well. The strength of the
231 model comes from the use of word embeddings. With numbers, that ability is neglected. Either the
232 data is changed to only take in strings or the model is equipped to process these numbers intelligently.

233 The next big issue is the specificity of the generated text. I need to come up with a way of injecting
234 named entities into the generated text. I might have to replace team and player entities with placeholder
235 tokens. This will probably be in the form of special substitute tokens that can be replaced with the
236 appropriate strings after the decoder stage.

237 The multiplicative attention of this model may not be the best mechanism to use for this task. It will
238 be interesting to see how different attention mechanisms affect the efficacy of the model.

239 The other question is, how do you introduce larger narratives into these summaries? There needs to
240 be a way for the NMT system to understand the importance of each and every moment to the larger
241 picture. That sounds difficult and it probably is. It is a key part of good sports-writing that probably
242 won't be mechanically reproducible for a while.

243 Sample output 3 really underscores how unreliable the training set maybe. My approach for creating
244 training examples was admittedly quite naive. There needs to be a better way of relating game events
245 with game summary sentences. It may require manual effort to do effectively. A simple regex does
246 not suffice.

247 **6.1 Other Datasets?**

248 An improved version will probably rely on data from other sources. While the Football Events Data
249 is comprehensive, it is not well suited to this task. Many event elements are incomplete - they have
250 'NA' in place of details. It is an inconsistent form that significantly decreases the effectiveness of
251 training.

252 Additionally, if the idea is to create a system that instantly generates summaries of football matches,
253 the Football Events set is probably not the most reliable. It appears to be manually derived and pretty
254 time intensive to create, which defeats the purpose of the system. Instead, the ideal system would
255 probably rely on data streams that are already in place. Using video input or voice commentary would
256 be interesting, though incredibly difficult. What is more feasible, and what I will probably explore
257 next, is using the live game feed - occasional text sentences that describe events - as source sentences
258 for a match summary NMT.

259 **7 Additional Information**

260 **Mentor:**

261 Michael Hermann Hahn - mhahn2@stanford.edu2

262 References

- 263 [1] Elisabeth Andre, Gerd Herzog, and Thomas Rist. “On the Simultaneous Interpretation of Real
264 World Image Sequences and their Natural Language Description: The System Soccer.” In:
265 *ECAI*. Vol. 88. 1988, pp. 449–54.
- 266 [2] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. “Neural machine translation by
267 jointly learning to align and translate”. In: *arXiv preprint arXiv:1409.0473* (2014).
- 268 [3] Kim Binsted. “A talking head architecture for entertainment and experimentation”. In: *AAAI*
269 *Fall Symposium Emotional and Intelligent: The Tangled Knot of Cognition*. 1998.
- 270 [4] Nadjat Bouayad-Agha, Gerard Casamayor, and Leo Wanner. “Content selection from an
271 ontology-based knowledge base for the generation of football summaries”. In: *Proceedings of*
272 *the 13th European Workshop on Natural Language Generation*. Association for Computational
273 Linguistics. 2011, pp. 72–81.
- 274 [5] Nadjat Bouayad-Agha et al. “Perspective-oriented generation of football match summaries:
275 Old tasks, new challenges”. In: *ACM Transactions on Speech and Language Processing (TSLP)*
276 9.2 (2012), p. 3.
- 277 [6] Kyunghyun Cho et al. “On the properties of neural machine translation: Encoder-decoder
278 approaches”. In: *arXiv preprint arXiv:1409.1259* (2014).
- 279 [7] Andreas Graefe. “Guide to automated journalism”. In: (2016).
- 280 [8] Melvin Johnson et al. “Google’s multilingual neural machine translation system: Enabling
281 zero-shot translation”. In: *Transactions of the Association for Computational Linguistics 5*
282 (2017), pp. 339–351.
- 283 [9] Nal Kalchbrenner and Phil Blunsom. “Recurrent continuous translation models”. In: *Proceed-*
284 *ings of the 2013 Conference on Empirical Methods in Natural Language Processing*. 2013,
285 pp. 1700–1709.
- 286 [10] Minh-Thang Luong, Hieu Pham, and Christopher D Manning. “Effective approaches to
287 attention-based neural machine translation”. In: *arXiv preprint arXiv:1508.04025* (2015).
- 288 [11] Kishore Papineni et al. “BLEU: a method for automatic evaluation of machine translation”.
289 In: *Proceedings of the 40th annual meeting on association for computational linguistics*.
290 Association for Computational Linguistics. 2002, pp. 311–318.
- 291 [12] Alin Secareanu. *Football Events*. [https://www.kaggle.com/secareanualin/football-](https://www.kaggle.com/secareanualin/football-events)
292 [events](https://www.kaggle.com/secareanualin/football-events). 2017.
- 293 [13] *SkySports Football Match Reports*.
- 294 [14] Ilya Sutskever, Oriol Vinyals, and Quoc V Le. “Sequence to sequence learning with neural
295 networks”. In: *Advances in neural information processing systems*. 2014, pp. 3104–3112.
- 296 [15] Yonghui Wu et al. “Google’s neural machine translation system: Bridging the gap between
297 human and machine translation”. In: *arXiv preprint arXiv:1609.08144* (2016).

298 Appendices

299 Football Event Legend

300	event_type	
301	0	Announcement
302	1	Attempt
303	2	Corner
304	3	Foul
305	4	Yellow card
306	5	Second yellow card
307	6	Red card
308	7	Substitution
309	8	Free kick won
310	9	Offside
311	10	Hand ball

312 11 Penalty conceded
 313
 314
 315 event_type2
 316 12 Key Pass
 317 13 Failed through ball
 318 14 Sending off
 319 15 Own goal
 320
 321
 322 side
 323 1 Home
 324 2 Away
 325
 326
 327 shot_place
 328 1 Bit too high
 329 2 Blocked
 330 3 Bottom left corner
 331 4 Bottom right corner
 332 5 Centre of the goal
 333 6 High and wide
 334 7 Hits the bar
 335 8 Misses to the left
 336 9 Misses to the right
 337 10 Too high
 338 11 Top centre of the goal
 339 12 Top left corner
 340 13 Top right corner
 341
 342
 343 shot_outcome
 344 1 On target
 345 2 Off target
 346 3 Blocked
 347 4 Hit the bar
 348
 349
 350 location
 351 1 Attacking half
 352 2 Defensive half
 353 3 Centre of the box
 354 4 Left wing
 355 5 Right wing
 356 6 Difficult angle and long range
 357 7 Difficult angle on the left
 358 8 Difficult angle on the right
 359 9 Left side of the box
 360 10 Left side of the six yard box
 361 11 Right side of the box
 362 12 Right side of the six yard box
 363 13 Very close range
 364 14 Penalty spot
 365 15 Outside the box
 366 16 Long range
 367 17 More than 35 yards
 368 18 More than 40 yards
 369 19 Not recorded
 370

```
371
372 bodypart
373 1      right foot
374 2      left foot
375 3      head
376
377
378 assist_method
379 0      None
380 1      Pass
381 2      Cross
382 3      Headed pass
383 4      Through ball
384
385
386 situation
387 1      Open play
388 2      Set piece
389 3      Corner
390 4      Free kick
```