

# ***Deep Learning* pentru descrierea automată a imaginilor în limbaj natural - *Image Captioning***

Anca Mihaela HOTĂRAN, Mihnea Horia VREJOIU

Institutul Național de Cercetare-Dezvoltare în Informatică – ICI București

anca.hotaran@ici.ro, mihnea.vrejoiu@ici.ro

**Rezumat:** Sintagma *Image Captioning* (IC) în contextul vederii artificiale se referă la generarea automată de descrieri textuale asociate imaginilor digitale. Nu este vorba doar de recunoașterea obiectelor din aceste imagini, ci și de descrierea atributelor lor, precum și a relațiilor și interacțiunilor dintre ele, totul exprimat textual în limbaj natural, corect din punct de vedere sintactic și semantic. Sintetic, pașii principali în generarea automată de descrieri textuale asociate imaginilor sunt: a) – extragerea informației vizuale din imagine și, b) – „traducerea” acesteia într-un text adecvat și semnificativ. Evoluțiile spectaculoase din domeniul rețelelor neuronale adânci și *Deep Learning* din ultimii ani au condus la progrese absolut remarcabile și în domeniul IC, calitatea textelor descriptive generate fiind îmbunătățită substanțial. Rețelele neuronale convoluționale (*Convolutional Neural Network* – CNN) au fost folosite în mod natural pentru obținerea de reprezentări vectoriale esențializate ale caracteristicilor (*features*) din imagini, iar rețelele neuronale recurente (*Recurrent Neural Network* – RNN), în particular de tip *Long Short-Term Memory* (LSTM), au fost utilizate pentru decodarea acestor reprezentări în fraze în limbaj natural. În lucrarea de față prezentăm o trecere în revistă a noilor tehnici și metode bazate pe *Deep Learning* utilizate în domeniul IC, cu detalierea și analizarea ca studiu de caz a uneia dintre cele mai performante dintre acestea, utilizând o arhitectură de tip *encoder-decoder* combinată cu un mecanism de focalizare a atenției vizuale pe regiunile corespunzătoare relevante din imagine la generarea fiecărui cuvânt nou din secvența de ieșire.

**Cuvinte cheie:** descriere textuală a imaginilor, învățare automată, învățare profundă, rețea neuronală adâncă, rețea convoluțională, rețea recurentă, LSTM, *encoder-decoder*, mecanism atențional.

## ***Deep Learning* for automatically describing images in natural language - *Image Captioning***

**Abstract:** *Image Captioning* (IC) in Computer Vision context refers to the automatic generation of textual descriptions associated with digital images. It is not only the recognition of the objects in these images, but also the description of their properties, as well as the relationships and interactions between them, all expressed textually in natural language, syntactically and semantically correct. Synthetically, the main steps in the automatic generation of textual descriptions associated with the images are: a) – extracting the visual information from the image, and, b) – “translating” it into an adequate and meaningful text. The spectacular developments in the field of deep neural networks and *Deep Learning* in recent years have led to absolutely remarkable progress also in the field of IC, the quality of the generated descriptive texts being substantially improved. *Convolutional Neural Networks* (CNN) have been naturally used to obtain essentialized vectorial representations of the image *features*, and *Recurrent Neural Networks* (RNN), in particular *Long Short-Term Memory* (LSTM), were used to decode these representations into phrases in natural language. In this paper we present an overview of the new techniques and methods based on *Deep Learning* used in the IC field, while also detailing and analyzing, as a case study, one of the best performing ones, using an *encoder-decoder* architecture combined with a mechanism for focusing the visual attention on the appropriate relevant regions of the image when generating each new word in the output sequence.

**Keywords:** image captioning, machine learning, deep learning, deep neural network, convolutional network, recurrent network, LSTM, encoder-decoder, attentional mechanism.

### **1. Introducere**

Generarea automată a unui text descriptiv asociat unei imagini (*Image Captioning* - IC), ca parte a obiectivului mai complex de înțelegere a scenelor în vederea artificială, este în sine o sarcină netrivială. Este vorba de mai mult decât de identificarea obiectelor care apar în imagini, respectiv de surprinderea și exprimarea atributelor lor, precum și a relațiilor și interacțiunilor dintre ele, textual în limbaj natural coerent, corect din punct de vedere sintactic și semantic.

Aria de aplicabilitate este destul de largă, pornind de la adnotarea descriptivă automată a volumului tot mai mare de imagini postate de utilizatorii individuali pe rețelele de *social media* sau crearea de rezumate textuale descriptive pentru bazele de date de imagini pentru organizare și/sau accesibilitate mai eficientă, asistarea persoanelor cu deficiențe vizuale majore prin descrierea textuală, combinată cu un sistem *text-to-speech* (TTS), a unor astfel de imagini, filme, sau chiar a mediului înconjurător în care se află sau se deplasează acestea, adnotarea cu descrieri ale celor observate în imagistica medicală de diferite tipuri pentru asistarea în punerea unor diagnostice și/sau urmărirea unor evoluții ș.a.m.d.

Înțelegerea conținutului unei imagini depinde în mod esențial de extragerea caracteristicilor specifice (*features*) din aceasta. Tehnicile de învățare automată (*Machine Learning* – ML) utilizate în acest scop pot fi încadrate în două mari categorii: tehnici bazate pe metode tradiționale și, respectiv, tehnici noi bazate pe *Deep Learning*. Nu vom aborda în această lucrare aspecte legate de metodele și tehnicile tradiționale, cum sunt de exemplu cele utilizând șabloane (*templates*) de descrieri – în care sunt introduse în poziții predefinite (*slot-uri*) denumirile obiectelor, atributele acestora sau relațiile între ele identificate în imaginea analizată, sau cele bazate pe căutarea și alegerea dintr-o mulțime predefinită a celei mai potrivite descrieri, asociate imaginii celei mai asemănătoare vizual cu cea analizată.

Problemele și etapele principale în generarea automată de descrieri textuale asociate imaginilor cu metode și tehnici bazate pe *Deep Learning* sunt: **a**) – extragerea informației vizuale din imagine și **b**) – transformarea acesteia într-un text adecvat și semnificativ. În ultimii ani, odată cu evoluțiile spectaculoase din domeniul rețelelor neuronale adânci și *Deep Learning* cu performanțe comparabile și chiar superioare celor umane în domenii ca analiza de imagini și traducerea automată, au fost înregistrate progrese absolut remarcabile și în domeniul IC, calitatea descrierilor textuale generate fiind îmbunătățită substanțial prin combinarea rețelelor neuronale convoluționale (*Convolutional Neural Network* – CNN), folosite pentru obținerea unor reprezentări vectoriale esențializate ale caracteristicilor (*features*) din imagini, cu rețele neuronale recurente (*Recurrent Neural Network* - RNN), în particular de tip *Long Short-Term Memory* (LSTM), utilizate pentru traducerea acestor reprezentări în fraze în limbaj natural.

Tehnicile de IC bazate pe *Deep Learning* pot fi grupate în mai multe categorii [Hossain et al., 2019], după: modul de mapare a caracteristicilor – *feature mapping* – (în spațiul vizual vs. într-un spațiu multimodal); tipul învățării (supervizată vs. nesupervizată); numărul de descrieri textuale și aria lor de acoperire sau extinderea (descrieri „dense” vs. descrieri ale întregii scene); arhitectura sistemului (*encoder-decoder* vs. arhitecturi compoziționale); modelul utilizat pentru partea de limbaj natural (RNN-LSTM vs. altele); alte abordări complementare utilizate (metode bazate pe atenție vs. metode bazate pe concepte semantice vs. extinderea scenariilor generate vs. metode de generare a unor descrieri textuale stilizate).

Lucrarea de față își propune o prezentare succintă a tehnicilor și metodelor bazate pe *Deep Learning* utilizate în domeniul IC, cu detalierea și analiza ca studiu de caz a uneia dintre cele mai performante dintre acestea, prima în care s-a utilizat și un mecanism atențional pe o arhitectură de tip *encoder-decoder*.

În continuare, lucrarea este structurată după cum urmează: Secțiunea 2 trece în revistă principalele aspecte privind noile metode și tehnici bazate pe *Deep Learning* utilizate în generarea automată de descrieri textuale ale imaginilor. În secțiunea 3 este analizată ca studiu de caz una dintre metodele de referință în IC [Xu et al., 2015], utilizând o arhitectură de tip *encoder-decoder* combinată cu un mecanism de focalizare dinamică a atenției pe regiunile corespunzătoare relevante ale imaginii la generarea fiecărui cuvânt din secvența descriptivă de ieșire. Lucrarea se încheie cu secțiunea 4 în care sunt enunțate câteva concluzii și comentarii.

## 2. Metode și tehnici de *Image Captioning* utilizând *Deep Learning*

Problema tratată de IC se referă în esență la extragerea dintr-o imagine de intrare a unei reprezentări a caracteristicilor esențiale și generarea automată pe baza acestei reprezentări a unei descrieri textuale a conținutului imaginii respective în limbaj natural. Aceasta se realizează în general ulterior unui proces corespunzător de învățare automată (*machine learning* – ML).

După modul de mapare a caracteristicilor, metodele de IC pot fi grupate în metode cu mapare în spațiul vizual sau într-un spațiu multimodal. În cazul mapării în spațiul vizual, caracteristicile imaginilor și descrierile textuale corespunzătoare sunt transmise independent *decoder*-ului de limbaj la învățare. În cazul spațiului multimodal, la antrenare, sistemul învață un spațiu comun obținut din imagini și descrierile asociate acestora, reprezentate în acest același spațiu prin ieșirile generate de un *encoder* de imagini și, respectiv, un *encoder* de limbaj, iar reprezentarea comună fiind transmisă *decoder*-ului.

Metodele de învățare automată utilizate în IC pot fi supervizate sau nesupervizate. În cazul învățării supervizate, datele de antrenare au atașate etichetele de clasă care reprezintă ieșirea așteptată, în timp ce, în cazul învățării nesupervizate, datele sunt neetichetate. Dintre metodele de învățare automată nesupervizată, cele mai folosite sunt *Generative Adversarial Networks* (GAN) și învățarea prin recompensă (*Reinforcement Learning* – RL). În această din urmă tehnică, învățarea este realizată de un agent prin explorarea datelor pas cu pas, alegând câte o acțiune în fiecare stare și primind „recompense” care evaluează acțiunea și/sau noua stare, obiectivul fiind maximizarea recompensei cumulate pe termen lung. Rețelele de tip GAN învață caracteristici din date neetichetate aplicând un mecanism competitiv între două rețele, generatorul și discriminatorul. Rețele neuronale bazate pe învățare supervizată au fost folosite cu succes în domenii precum clasificarea imaginilor, detectarea de obiecte sau învățarea atributelor, obținând performanțe comparabile cu, sau chiar peste cele umane. Această experiență a fost aplicată și în domeniul generării de descrieri textuale asociate imaginilor. Arhitecturile utilizate pot fi: de tip *encoder-decoder*; compoziționale; cu mecanisme bazate pe atenție; bazate pe concepte semantice; cu metode de generare a unor descrieri textuale stilizate; cu extinderea scenariilor generate; cu generare de descrieri „dense”.

Descrierile textuale ale imaginilor pot fi „dense” sau globale ale întregii scene. În cazul descrierii textuale dense sunt localizate regiunile importante ale imaginii și se generează descrieri pentru fiecare dintre acestea. Pașii urmați sunt în general următorii: se generează propuneri de regiuni de interes (*regions of interest* – ROI) din imagini; o rețea neuronală convoluțională extrage caracteristicile din regiunile selectate; ieșirile generate sunt utilizate de un model al limbajului pentru a genera descrieri textuale pentru fiecare regiune. Descrierea globală a întregii imagini poate să fie subiectivă și să nu ofere suficientă informație pentru înțelegerea acesteia. Descrierile asociate regiunilor de interes sunt mai obiective și mai detaliate. Există desigur anumite dificultăți și provocări în generarea descrierilor dense. Printre acestea, menționăm faptul că un obiect poate avea mai multe regiuni de interes care se suprapun, sau faptul că este dificil de recunoscut fiecare regiune pentru toate conceptele vizuale. În cazul generării de descrieri textuale corespunzătoare întregii imagini, metodele cele mai utilizate sunt: arhitecturi de tip *encoder-decoder*; arhitecturi compoziționale; arhitecturi bazate pe semantică; arhitecturi bazate pe atenție; descrieri textuale stilizate; extinderea scenariilor.

Cele mai multe modele se bazează pe arhitectura *encoder-decoder*, care este foarte flexibilă și performantă. Această arhitectură a fost inițial proiectată pentru traducerea automată a limbajului. În cazul IC putem vorbi de o traducere automată a unei imagini într-un text descriptiv corespunzător. Astfel, generarea automată a descrierilor textuale ale imaginilor necesită logic două modele: un model pentru extragerea caracteristicilor din imagini (*encoder*) și un model al limbajului, care translatează caracteristicile și obiectele furnizate de modelul de tratare a imaginilor într-o frază în limbaj natural (*decoder*). *Encoder*-ul este, cel mai adesea, o rețea neuronală convoluțională pentru extragerea caracteristicilor din imagini, detectarea obiectelor și a relațiilor dintre acestea. În lucrarea [Staniute & Sesok, 2019] este prezentată o analiză sistematică a literaturii de specialitate din domeniul IC publicate în ultimii patru ani. Sunt evidențiate cele mai folosite tehnici și cele mai importante provocări în domeniul IC. Au fost identificate cinci tipuri de rețele convoluționale utilizate pentru extragerea de caracteristici din imagini, dintre care două majoritare: VGGNet și ResNet. Acestea au fost raportate ca utilizate fiecare în câte 33 de articole din cele 78 analizate, primul fiind preferat pentru simplitatea modelului și pentru puterea sa, iar al doilea pentru numărul mai mic de parametri și eficiența sa computațională în comparație cu toate celelalte arhitecturi convoluționale. *Decoder*-ul este în general o rețea neuronală recurentă

(*Recurrent Neural Network* - RNN), care permite reprezentarea contextului anterior, memorând secvența de cuvinte anterioare pentru a prezice următorul cuvânt din secvență. Arhitecturi particulare de RNN, precum *Long Short-Term Memory* (LSTM), rezolvă problema diminuării excesive a gradientului (*vanishing gradient*) într-un lanț lung de transformări neliniare, învățând astfel dependențe pe termen lung în secvența de date. Utilizarea LSTM este considerată cea mai populară metodă în contextul IC. Cu toate că rețelele de tip LSTM posedă memorie în care se poate păstra informație despre secvența deja generată mai bine decât în cazul rețelele convoluționale, este însă necesară actualizarea la fiecare pas, fiind un model consumator de memorie. De aceea au fost reconsiderate unele sisteme care utilizau combinația CNN-RNN, fiind propuse în locul acestora arhitecturi CNN-CNN. Rețelele convoluționale pot fi antrenate mult mai rapid decât cele recurente, în principal din două motive: convoluțiile pot fi procesate în paralel, pe când recurențele numai secvențial; procesoarele grafice (GPU) pot fi folosite pentru a crește viteza de antrenare a modelelor convoluționale, în timp ce (deocamdată) nu există *hardware* care să accelereze antrenarea rețelelor recurente. Totuși, rezultatele nu au fost cele dorite. Din cele 78 de sisteme considerate, 68 au revenit la LSTM. Primele sisteme de tip *encoder-decoder* au utilizat LSTM unidirecționale pentru generarea textului, cuvântul următor fiind prezis pe baza contextului vizual și cuvintelor generate anterior. În [Wang et al., 2016] este propus un model bidirecțional cu două rețele LSTM separate, care utilizează atât informația contextuală trecută cât și cea viitoare pentru generarea descrierii textuale. Această metodă permite generarea unor descrieri mai bogate, atât semantic cât și contextual. Viteza depinde nu numai de caracteristicile modelului, ci și de dimensiunea vocabularului. În [Mishra & Liwicki, 2019] s-a încercat reducerea dicționarului de cuvinte utilizate de la 10.000–40.000 de cuvinte la numai 248 de cuvinte. Deși este o reducere drastică, rezultatele au fost încurajatoare.

Arhitecturile compoziționale sunt alcătuite din mai multe blocuri funcționale independente. Pașii urmați sunt: o rețea neuronală convoluțională extrage caracteristicile din imagine; conceptele vizuale (atribute) se obțin din caracteristici; se generează un set de descrieri textuale utilizând informațiile din pașii anteriori; descrierile sunt reevaluate utilizând un model de similaritate multimodal, care selectează descrierile cele mai adecvate.

Utilizarea mecanismelor de atenție vizuală a permis obținerea de performanțe superioare. Acestea au la bază comportamentul uman natural: oamenii acordă o atenție specială anumitor regiuni ale imaginii și formulează explicații ale relațiilor dintre obiectele ce apar în aceste regiuni. Astfel, mecanismele bazate pe atenție au devenit foarte populare în *Deep Learning* pentru că permit depășirea unor limitări prin focalizarea dinamică asupra diverselor regiuni din imaginea de intrare pe parcursul generării secvenței de cuvinte la ieșire. Pașii urmați sunt: o rețea neuronală convoluțională extrage caracteristicile din întreaga imagine; modelul limbajului generează cuvinte și fraze pornind de la aceste caracteristici; cu fiecare cuvânt generat, mecanismul de atenție determină regiunea de interes din imagine corespunzătoare descrierii curente pentru generarea cuvântului următor; această regiune, împreună cu textul deja generat, sunt utilizate de modelul limbajului în generarea cuvântului următor. Majoritatea modelelor cu mecanisme atenționale pentru generarea descrierilor textuale asociate imaginilor utilizează imaginea la fiecare pas, chiar și atunci când se generează cuvinte de legătură, exprimând relații ș.a. care nu au un corespondent vizual în imagine. În [Lu et al., 2016] se propune introducerea unei „santinele” vizuale, care să stabilească dacă la un anumit moment modelul se bazează pe semnalul vizual sau pe modelul limbajului, implementând un model de atenție adaptivă. Sistemul extinde modelul LSTM al *decoder*-ului cu o nouă poartă santinelă, care decide ponderea informației extrase de *decoder* din imagine pentru a genera cuvântul următor. Mecanismele atenționale pot fi de două tipuri: stocastic „*hard*” sau determinist „*soft*”. Pe de altă parte, abordările atenționale pot fi de tip *top-down* sau *bottom-up*. [He et al., 2019] au confirmat recent că abordarea *top-down* este mai eficientă, rezultatele fiind similare cu cele din experimentele cu subiecți umani. Sistemele de IC folosesc în general mecanisme atenționale de tip *top-down*. În [Anderson et al., 2018] este descris un sistem care propune un mecanism atențional ce combină metodele *top-down* și *bottom-up* pentru stabilirea regiunilor de interes la nivel de obiecte sau alte zone ale imaginii. Mecanismul *bottom-up* utilizează o rețea convoluțională bazată pe regiuni (*region based*) optimizată – Faster R-CNN [Ren et al., 2015] – pentru a propune împărțirea imaginii în regiuni de interes și generarea vectorilor de

caracteristici asociată fiecăreia dintre acestea, în timp ce mecanismul *top-down* utilizează LSTM pentru determinarea ponderilor caracteristicilor. În fine, nivelul limbajului este implementat la rândul său utilizând un al doilea LSTM.

În ultimul an, doi termeni apar în aproape fiecare articol din domeniul IC: extinderea scenariilor generate (*novel*) și semantică [Hossain et al., 2019]. Acești termeni au apărut în legătură cu provocarea actuală cea mai importantă: generarea automată de descrieri textuale ale imaginilor care să fie cât mai apropiate de cele produse de subiecți umani. Implementarea semantică facilitează o modalitate de a introduce sentimente în sistemul de generare automată de descrieri textuale. Extinderea scenariilor generate încearcă depășirea limitărilor sistemelor actuale, datorate mai multor motive: în primul rând, modelele sunt antrenate pe seturi de date specifice, care nu acoperă toate scenariile posibile și, de asemenea, dicționarele de cuvinte sunt limitate ca număr de cuvinte și combinații; în al doilea rând, modelele sunt în general gândite să execute o singură sarcină specifică, în timp ce oamenii sunt capabili să efectueze mai multe sarcini simultan. Au fost propuse abordări pentru acest aspect al IC, care în principal propun un dialog cu utilizatorul, prin care sistemul pune întrebări specifice și învață din răspunsuri. Această metodă este specifică domeniului *question-answering*, dar are un impact foarte însemnat pentru generarea automată de descrieri textuale asociate imaginilor. Cu toate că metodele de IC bazate pe *Deep Learning* au avut rezultate încurajatoare, acestea depind în mare măsură de seturile de date imagine-descrieri. Aceste metode pot genera descrieri numai pentru obiectele care au apărut în setul de date de antrenare și, ca urmare, necesită seturi foarte mari de date. Metodele care extind contextul scenariilor pot genera descrieri ale obiectelor ce nu apar în setul de antrenare. În general, pașii urmați sunt: se antrenează separat un clasificator lexical și un model al limbajului pe seturi de imagini și texte necorelate; se antrenează un model de generare de descrieri textuale pe seturi de date imagine-descrieri; în final, ambele modele se combină pentru a fi antrenate împreună. Generarea de descrieri textuale pentru obiecte care nu au mai fost văzute este un subiect de cercetare recent. În lucrarea [Venugopalan et al., 2017] este introdus *Novel Object Captioner* (NOC) pentru generarea de descrieri ale obiectelor noi utilizând surse externe pentru recunoașterea acestor obiecte și folosind cunoaștere semantică. Metodele bazate pe concepte semantice folosesc selectiv un set de propuneri de concepte semantice extrase din analiza imaginilor. Aceste concepte sunt apoi atașate stărilor ascunse și ieșirilor rețelei neuronale recurente. Etapele acestor metode sunt: o rețea neuronală convoluțională *encoder* extrage caracteristicile imaginilor și conceptele semantice; aceste caracteristici sunt transmise ca intrări unui model de generare a limbajului; conceptele semantice sunt adăugate unor stări ascunse ale modelului limbajului; rețeaua responsabilă cu generarea limbajului produce descrierile textuale cu concepte semantice.

Majoritatea sistemelor de IC se bazează strict pe caracteristicile extrase din imagini, rezultând o descriere factuală a acestora. Descrierile stilizate pot fi însă mult mai expresive și atractive decât o descriere factuală. Astfel de descrieri sunt apreciate în multe aplicații din lumea reală. De exemplu, oamenii încarcă o mulțime de fotografii în diverse sisteme *social media*. Aceste imagini necesită descrieri stilizate și atractive. Aceste metode folosesc, în general, următorii pași: o rețea convoluțională extrage caracteristicile imaginii; este folosit un corpus separat de text pentru a extrage diverse concepte stilistice (de ex. stil romantic, umoristic ș.a.) din datele de antrenare; din informațiile de la pașii anteriori modelul limbajului poate genera descrieri textuale stilizate.

Cele mai multe sisteme s-au limitat la a folosi exclusiv seturile de date specifice pentru IC, restrângând astfel categoriile de obiecte pentru care pot fi generate descrieri textuale. Metoda propusă de [Venugopalan et al., 2017] poate utiliza suplimentar surse de date externe, precum imagini etichetate din seturi de date folosite la clasificarea/recunoașterea de obiecte, sau informații semantice extrase din texte neadnotate.

### **3. Image Captioning utilizând o arhitectură *encoder-decoder* cu mecanism atențional**

În cele ce urmează prezentăm ca studiu de caz o metodă de descriere automată a imaginilor (*Image Captioning* – IC) având la bază ideile și rezultatele prezentate de [Xu et al., 2015]. Este

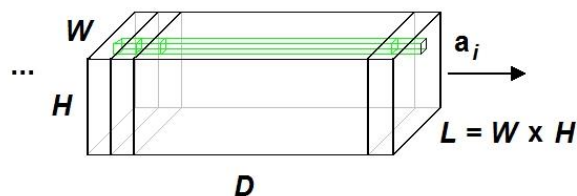


vorba despre o abordare *Deep Learning encoder-decoder*, îmbunătățită însă prin introducerea unui **mecanism atențional** care permite focalizarea dinamică pe o anumită regiune a reprezentării imaginii de intrare pentru predicția următorului cuvânt din secvența descriptivă de ieșire.

Arhitectura de tip *encoder-decoder* a fost inițial utilizată pentru învățarea de tip secvență la secvență (*sequence to sequence*) cu rețele neuronale recurente utilizată în traducerea automată [Cho et al., 2014; Bahdanau et al., 2014; Sutskever et al., 2014]. Generarea descrierii unei imagini este privită analog unei „traduceri” a imaginii într-o frază descriptivă a acesteia. Astfel, la intrarea modelului (a *encoder*-ului) este aplicată imaginea  $x$ , iar la ieșirea acestuia (a *decoder*-ului) este generată o secvență de cuvinte în limbaj natural asociate respectivei imagini:  $y = \{y_1, y_2, \dots, y_C\}$ ,  $y_i \in \mathbf{R}^K$ , unde  $K$  este dimensiunea vocabularului de cuvinte disponibile, iar  $C$  este lungimea secvenței de cuvinte care descrie imaginea.

Metode anterioare (de ex. [Vinyals et. al, 2014]) s-au bazat în general pe o unică reprezentare statică comprimată, obținută în stratul final complet conectat al unei rețele convoluționale [Vrejoiu, 2019a] utilizate ca *encoder*, care codifică imaginea de intrare integrală și este utilizată ca intrare de *decoder* doar o singură dată, la primul pas. Totuși, astfel se poate pierde informație potențial utilă pentru descrieri mai cuprinzătoare. Prin utilizarea unor reprezentări dintr-un nivel convoluțional precedent poate fi valorificată inclusiv astfel de informație, fiind necesară însă și „direcționarea” modelului către ceea ce este important în acestea la generarea fiecărui nou cuvânt la fiecare pas. Inspirată din funcționarea sistemului vizual uman [Rensink, 2000; Corbetta & Shulman, 2001], focalizarea selectivă prin mecanisme atenționale permite caracteristicilor semnificative dintr-o imagine să fie puse în evidență în mod dinamic contextual. Implementarea analizată de noi utilizează un mecanism atențional determinist de tip „soft” (introdus de [Bahdanau et al., 2014]), antrenabil prin metode standard folosind retropropagarea erorii (*backpropagation*).

**Partea de *encoder*** este reprezentată de o **rețea neuronală convoluțională** (*Convolutional Neural Network – CNN*) fără straturile finale de agregare și complet conectate, aceasta terminându-se cu o hartă de caracteristici (*feature map*) cu  $D$  „canale” (straturi de neuroni), de dimensiuni spațiale  $W \times H$ , corespunzătoare aplicării ultimilor  $D$  filtre de convoluție. Pentru fiecare imagine de intrare, este generat la ieșire un set de  $L = W \times H$  vectori numiți **vectori de adnotare** (*annotation vectors*), fiecare de dimensiune  $D$  și reprezentând câte o porțiune din respectiva imagine, localizată în poziția  $i$ :  $a = \{a_1, a_2, \dots, a_L\}$ ,  $a_i \in \mathbf{R}^D$ . Fiecare componentă  $a_i[j]$ , ( $j = 1 \div D$ ) a unui vector  $a_i$  reprezintă un indicator pentru detectarea (în mai mare sau mai mică măsură) a caracteristicii  $j$  în locația  $i$  ( $i = 1 \div L$ ) din imagine. În figura 1 este reprezentată intuitiv structura ultimului nivel convoluțional considerat (harta de caracteristici – *feature map*) de dimensiuni  $L \times D$ , ( $L = W \times H$ ) și ieșirea *encoder*-ului (vectorii de adnotare  $a_i$ ), fiecare secțiune transversală corespunzând unei caracteristici  $j$ , ( $j = 1 \div D$ ).



**Figura 1.** Ultimul nivel (convoluțional) și ieșirea *encoder*-ului

Această ieșire a *encoder*-ului este utilizată mai departe de către *decoder*. Prin utilizarea ieșirii unui strat convoluțional anterior în locul unui final complet conectat, se asigură posibilitatea *decoder*-ului de a se putea focaliza dinamic selectiv pe anumite porțiuni din imagine prin alegerea unui subset din toți vectorii de adnotare.

**Partea de *decoder*** este reprezentată de o **rețea neuronală recurentă** (*Recurrent Neural Network – RNN*) de tip **LSTM** (*Long Short-Term Memory*) [Hochreiter & Schmidhuber, 1997; Vrejoiu, 2019b], care pentru fiecare reprezentare de imagine aplicată la intrare produce o descriere textuală. La fiecare pas de timp  $t$  este generat câte un cuvânt,  $y_t$ , pe baza cuvântului generat la pasul anterior,  $y_{t-1}$ , a stării ascunse (*hidden state*) obținute la pasul anterior,  $h_{t-1}$ , precum și a unui **vector de context** vizual,  $\hat{z}_t$ , care surprinde informația extrasă din imagine relevantă pentru generarea

noului cuvânt. Trebuie precizat faptul că pentru cuvinte este utilizată o reprezentare codificată numeric (*embedding*)  $\mathbf{E}y_{t-1}$ , respectiv  $\mathbf{E}y_t$ . Toate aceste trei mărimi vectoriale sunt aplicate ca intrare, modulate în fiecare caz prin ponderi antrenabile, la toate porțile celulei LSTM: de intrare, de uitare/ștergere și de ieșire. Astfel, celula LSTM învață cum să pondereze fiecare componentă de intrare (*input gate*), învățând totodată cum să moduleze respectiva contribuție la memorie (*input modulator*). De asemenea sunt învățate ponderi pentru ștergerea memoriei celulei LSTM (*forget gate*) și ponderi pentru controlul modului în care această memorie trebuie să emită (*output gate*). Funcționarea LSTM este descrisă la fiecare pas de timp  $t$  după cum urmează:

$$\mathbf{i}_t = \sigma(Wy_i \mathbf{E}y_{t-1} + Wh_i \mathbf{h}_{t-1} + Wz_i \hat{\mathbf{z}}_t + \mathbf{b}_i) \quad (1)$$

$$\mathbf{f}_t = \sigma(Wy_f \mathbf{E}y_{t-1} + Wh_f \mathbf{h}_{t-1} + Wz_f \hat{\mathbf{z}}_t + \mathbf{b}_f) \quad (2)$$

$$\mathbf{o}_t = \sigma(Wy_o \mathbf{E}y_{t-1} + Wh_o \mathbf{h}_{t-1} + Wz_o \hat{\mathbf{z}}_t + \mathbf{b}_o) \quad (3)$$

$$\mathbf{g}_t = \tanh(Wy_c \mathbf{E}y_{t-1} + Wh_c \mathbf{h}_{t-1} + Wz_c \hat{\mathbf{z}}_t + \mathbf{b}_c) \quad (4)$$

unde:  $Wy_{m,n}$ ,  $Wh_{m,n}$  și  $Wz_{D,n}$  reprezintă matricile de ponderi antrenabile, iar  $\mathbf{b}$  *bias*-urile globale asociate fiecărei porți ( $i, f, o$ ), respectiv memoriei celulei LSTM ( $c$ ), iar:

$$\mathbf{c}_t = \mathbf{f}_t \cdot \mathbf{c}_{t-1} + \mathbf{i}_t \cdot \mathbf{g}_t \quad (5)$$

$$\mathbf{h}_t = \mathbf{o}_t \cdot \tanh(\mathbf{c}_t) \quad (6)$$

Vectorii  $\mathbf{i}_t$ ,  $\mathbf{f}_t$ ,  $\mathbf{o}_t$ ,  $\mathbf{c}_t$ ,  $\mathbf{h}_t \in \mathbf{R}^n$ , de dimensiune  $n$  fiecare, reprezintă „semnalul” prin porțile de intrare, uitare/ștergere și ieșire, respectiv starea memoriei și starea ascunsă pentru rețeaua LSTM, iar  $\mathbf{g}_t \in \mathbf{R}^n$  modulează poarta de intrare la momentul  $t$ . Vectorul  $\hat{\mathbf{z}}_t \in \mathbf{R}^D$  este vectorul de context conținând informația vizuală asociată unei locații particulare din intrare la fiecare moment  $t$ .  $\mathbf{E} \in \mathbf{R}^{m \times K}$  reprezintă matricea de *embedding* corespunzătoare vocabularului de  $K$  cuvinte disponibile,  $m$  și  $n$  reprezintă dimensiunea de *embedding* și respectiv dimensiunea LSTM, iar  $\sigma$  și  $\tanh$  reprezintă funcțiile de activare de tip sigmoidă și respectiv tangentă hiperbolică și  $\cdot$  semnifică înmulțirea element cu element.

Vectorul de context  $\hat{\mathbf{z}}_t$  oferă o reprezentare dinamică a porțiunii relevante din imaginea de intrare pentru fiecare moment  $t$ . Acesta se calculează utilizându-se un mecanism  $\Phi$  pe baza vectorilor de adnotare  $\mathbf{a}_i$ , ( $i = 1 \div L$ ) corespunzători caracteristicilor extrase pentru diferite locații  $i$  din imagine. Pentru fiecare locație, este generată o pondere pozitivă, subunitară,  $\alpha_{ti}$ , care – în cazul mecanismului atențional determinist „soft” utilizat – indică importanța relativă ce trebuie acordată locației  $i$  la momentul  $t$ . Ponderea  $\alpha_{ti}$  a fiecărui vector de adnotare  $\mathbf{a}_i$  este calculată utilizându-se un **model atențional**  $f_{att}$  reprezentat de un perceptron multistrat (*Multi-Layer Perceptron* – MLP), pe baza stării ascunse anterioare a rețelei LSTM,  $\mathbf{h}_{t-1}$ . Această stare evoluează pe măsură ce rețeaua recurentă avansează în secvența de cuvinte de ieșire, locația unde rețeaua trebuie să caute mai departe depinzând astfel de secvența de cuvinte deja generate.

$$e_{ti} = f_{att}(\mathbf{a}_i, \mathbf{h}_{t-1}) = W_{att} ReLU(T_a \mathbf{a}_i + T_h \mathbf{h}_{t-1}), \quad (7)$$

unde: *ReLU* este funcția de activare de tip *Rectified Linear Unit*,  $T_a$  și  $T_h$  reprezintă ponderile a două straturi liniare separate asociate transformărilor aplicate fiecărui vector de adnotare  $\mathbf{a}_i$  de dimensiune  $D$  și, respectiv, stării ascunse anterioare a LSTM,  $\mathbf{h}_{t-1}$ , de dimensiune  $n$  pentru a fi combinate în spațiul (dimensionalitatea) modelului atențional, iar  $W_{att}$  reprezintă ponderile stratului final, cu câte 1 neuron pentru fiecare pereche de intrare ( $\mathbf{a}_i, \mathbf{h}_{t-1}$ ), care furnizează la ieșire valoarea  $e_{ti}$ , ponderi care sunt comune pentru fiecare locație spațială  $i$ . Ponderile  $\alpha_{ti}$  sunt obținute prin aplicarea funcției *softmax* peste cele  $L$  ieșiri  $e_{ti}$  ale modelului atențional corespunzătoare aplicării celor  $L$  vectori de adnotare  $\mathbf{a}_i$  (fiecare în conjuncție cu starea anterioară  $\mathbf{h}_{t-1}$ ), la intrarea acestuia:

$$\alpha_{ti} = \exp(e_{ti}) / \sum_{k \leq L} \exp(e_{tk}). \quad (8)$$

Cu ponderile  $\alpha_{ti}$  astfel calculate ( $0 \leq \alpha_{ti} \leq 1$  și  $\sum_{i \leq L} \alpha_{ti} = 1$ ), poate fi obținut vectorul de context la pasul  $t$ ,  $\hat{\mathbf{z}}_t$ , care „integrează” contribuțiile ponderate ale tuturor vectorilor de adnotare. Acesta nu mai depinde de locație, iar valorile celor  $D$  componente ale sale pun în evidență acele caracteristici/canale (reprezentând eventual obiecte) care sunt cele mai reprezentative contextual





$$L_d = -\log(P(\mathbf{y}|\mathbf{x})) + \lambda \sum_{i \leq L} (1 - \sum_{t \leq C} \alpha_{it})^2. \quad (15)$$

Se utilizează metoda gradientului descendent stocastică (*stochastic gradient descent* – SGD) și algoritmi cu rată de învățare adaptivă și optimizare RMSProp [Tieleman & Hinton, 2012] sau Adam (*Adaptive Moment Estimation*) [Kingma & Ba, 2014], precum și regularizare de tip *dropout* [Srivastava et al., 2014].

### Implementarea experimentală testată

Pentru experimentarea metodei prezentate mai sus s-a utilizat o implementare PyTorch [Vinodababu, 2018] dezvoltată pe baza celei inițiale în Theano a [Xu et al., 2015] și un sistem Dual Xeon 2,4GHz, cu 128GB DDR4 și placă GPU nVidia GeForce GTX 1080 Ti 11GB. A fost utilizat setul de date *Microsoft Common Objects in COntext* (MSCOCO), cu 91 de categorii de obiecte în 328.000 de imagini adnotate cu un total 2.500.000 de descrieri. Au fost folosite efectiv 123.287 imagini, 82.783 de antrenare și 40.504 de validare și testare, cu câte 5 descrieri textuale fiecare.

Ca și în cazul VGG și pentru ResNet imaginile de intrare de rezoluții variate sunt inițial redimensionate astfel încât latura cea mai mică să aibă 256 pixeli, cu păstrarea factorului de formă (*aspect ratio*). Intrarea *encoder*-ului este dată de cele 3 canale RGB obținute din decuparea regiunii centrale de  $224 \times 224$  pixeli a acestei imagini. Valorile pixelilor sunt aduse în intervalul de valori reale  $[0, 1]$  și apoi normalizate pe baza valorii medii și a deviației standard pentru fiecare canal (care se obțin utilizându-se tot setul de imagini de care se dispune). Astfel, PyTorch folosind convenția NCHW (N = dimensiunea *mini-batch*-ului, C = numărul de canale, H = înălțimea și W = lățimea imaginii), la intrarea modelului sunt aplicați tensori de dimensiuni N, 3, 256, 256, conținând valori reale din intervalul  $[0, 1]$ .

Pentru partea de *encoder* care furnizează vectorii de adnotare  $\mathbf{a}_i$  utilizați de *decoder* s-a folosit arhitectura convoluțională ResNet-101 [He et al., 2015], câștigătoare a competiției *ImageNet Large Scale Visual Recognition Challenge* (ILSVRC) 2015, mai performantă decât VGG-19 [Simonyan & Zisserman, 2014] utilizată de [Xu et al., 2015]. În experimente s-a utilizat practic implementarea PyTorch din modulul *torchvision* a modelului ResNet-101, preantrenat pe setul de date ImageNet 2014, fără ultimele două niveluri (de *average pooling* și liniar final) deoarece se dorește doar codificarea imaginilor și nu clasificarea acestora. *Encoder*-ul poate fi de asemenea antrenat de la început, sau doar ajustat fin (*fine-tuned*) odată cu restul modelului. În cel din urmă caz, trebuie ajustate fin numai ponderile pentru ultimele trei niveluri convoluționale „*bottleneck*” ale ResNet-101, deoarece în primele niveluri s-au învățat practic numai caracteristici simple fundamentale (linii, muchii, curbe etc.), care nu mai trebuie afectate în vreun fel. În cazul ResNet-101, numărul de canale din ultimul strat al ultimului nivel convoluțional „*bottleneck*” (Conv 5), care reprezintă ieșirea *encoder*-ului, este  $D = 2048$ , la o rezoluție spațială  $W \times H = 7 \times 7$ , față de  $D = 512$  și  $W \times H = 14 \times 14$  la VGG-19. Astfel, pentru a păstra o rezoluție analogă, se aplică o redimensionare adaptivă asupra acestei ieșiri la o rezoluție impusă de  $14 \times 14$ . Pentru creșterea vitezei, performanței și stabilității, fiecare strat convoluțional din *encoder* este urmat de aplicarea unui strat de normalizare pe lot (*batch normalization*).

Șirurile descriptive asociate imaginilor (*captions*) constituie atât intrări cât și ieșiri țintă pentru *decoder*, fiecare cuvânt fiind utilizat pentru predicția celui următor. Primul cuvânt generat este întotdeauna unul standard care marchează începutul descrierii,  $y_0 = \langle \text{start} \rangle$ , iar finalul descrierii este marcat întotdeauna prin generarea unui cuvânt standard corespunzător,  $y_{C+1} = \langle \text{end} \rangle$ , deoarece *decoder*-ul trebuie să învețe să prezică sfârșitul unui șir descriptiv (*caption*) pentru a se putea opri procesul inferențial de decodare odată prezis acesta:  $\langle \text{start} \rangle y_1 y_2 \dots y_C \langle \text{end} \rangle$ . Totodată, deoarece în PyTorch șirurile descriptive sunt manipulate ca tensori de dimensiuni fixe, este necesar ca acestea să fie completate, când este necesar, până la o aceeași lungime maximă stabilită,  $C_{max}$ , cu un cuvânt standard dedicat,  $y_j = \langle \text{pad} \rangle$ , ( $j = C+2 \div C_{max}$ ), astfel:  $\langle \text{start} \rangle y_1 y_2 \dots y_C \langle \text{end} \rangle \langle \text{pad} \rangle \langle \text{pad} \rangle \dots \langle \text{pad} \rangle$ .

Pe de altă parte, se crează o tabelă asociativă de indexare a cuvintelor din vocabularul utilizat, prin care fiecăruia dintre acestea, inclusiv celor trei standard:  $\langle \text{start} \rangle$ ,  $\langle \text{end} \rangle$  și  $\langle \text{pad} \rangle$ , le

este asociat câte un indice în această tabelă. Utilizând acest index prin care este astfel codat asociativ fiecare cuvânt prin indicele respectiv, este posibilă și asocierea *embedding*-ului corespunzător fiecăruia și identificarea locului ocupat în scorurile cuvintelor prezise.

Astfel, descrierile textuale ale imaginilor (*captions*) sunt tensori de întregi de dimensiuni  $N_c$ ,  $L_c$ , unde  $L_c$  este lungimea completată cu <pad>-uri până la sfârșit. Lungimea șirurilor descriptive (*caption lengths*) inclusiv cuvintele standard de început și sfârșit ( $2 +$  lungimea efectivă  $C$ ) este de asemenea menținută într-un tensor de întregi de dimensiune  $N_c$ . Astfel este posibilă procesarea unei secvențe doar până la această lungime fără a se mai pierde timp de calcul inutil pentru cuvintele de completare <pad> și de asemenea pot fi construite grafuri dinamice cu PyTorch.

Modelul atențional este unul simplu, compus doar din straturi liniare și câteva activări. Astfel, două straturi liniare separate transformă fiecare din cei  $14 \times 14 = 196$  vectori de adnotare  $\mathbf{a}_i$ , cu câte 2048 componente fiecare, obținuți la ieșirea *encoder*-ului și respectiv starea ascunsă (*hidden state*) obținută la ieșirea *decoder*-ului la momentul anterior,  $\mathbf{h}_{t-1}$ , de dimensiune 512, la o aceeași dimensiune atențională, 512. Rezultatele acestor transformări liniare sunt adunate și apoi rezultatul sumării este trecut printr-o activare de tip ReLU. În fine, un al treilea strat liniar transformă acest din urmă rezultat la dimensiunea  $1 - e_{ii}$  din relația (7) – pentru fiecare pereche ( $\mathbf{a}_i$ ,  $\mathbf{h}_{t-1}$ ). În fine, asupra tuturor acestor  $L$  valori  $e_{ii}$ , de dimensiune 1, se aplică funcția *softmax* pentru a fi generate ponderile  $\alpha_{ii}$  – date de relația (8), pentru fiecare locație spațială  $i$  la pasul  $t$ .

*Decoder*-ul primește la intrare ieșirea *encoder*-ului reprezentată de imaginea codificată în vectorii de adnotare  $\mathbf{a}_i$ , care este liniarizată la dimensiunile  $N$ , 196, 2048. Inițializarea stării ascunse,  $\mathbf{h}_0$  și a memoriei celulei LSTM,  $\mathbf{c}_0$ , se realizează, prin două straturi liniare separate, pe baza mediei vectorilor de adnotare  $\mathbf{a}_i$  furnizați la ieșirea *encoder*-ului – relațiile (11) și (12).

Deoarece implementarea necesită la învățare durate proporționale cu lungimea descrierii celei mai lungi la fiecare actualizare a ponderilor, antrenarea pe grupuri aleatorii de descrieri este evident ineficientă ca și consum de timp. Pentru a se evita aceasta, s-a efectuat o preprocesare prin care s-a construit mai întâi o tabelă (un dicționar) care asociază o anumită lungime subsetului corespunzător de descrieri de acea lungime. Apoi, la antrenare se eșantionează aleator o anumită astfel de lungime și se compun mini loturi (*mini-batches*) pe baza respectivei lungimi la fiecare pas, ceea ce îmbunătățește semnificativ viteza de convergență fără a afecta negativ performanța. Cele  $N$  imagini și descrierile (*captions*) asociate din fiecare lot (*batch*) sunt ordonate descrescător pe baza lungimii descrierilor respective, astfel încât să poată fi efectuate procesări numai pentru pașii de timp utili, respectiv să nu fie procesate și eventualele cuvinte de completare <pad>. Astfel, se poate itera la fiecare pas de timp numai pentru lungimea efectivă considerată la momentul respectiv, ceea ce determină dimensiunea efectivă  $N_t$  a lotului la acel pas. Ordonarea permite ca primele  $N_t$  imagini la fiecare pas să fie aliniată cu ieșirile de la pasul anterior (primele  $N_t$  imagini sunt procesate utilizându-se primele  $N_t$  ieșiri de la pasul anterior). Această iterație este efectuată explicit pas cu pas, individual, într-o buclă *for* aplicându-se mecanismul atențional între fiecare doi pași consecutivi de decodare.

Ponderile  $\alpha_{ii}$  și codificarea ponderată atențional la fiecare pas de timp sunt calculate de rețeaua atențională. Vectorul de context  $\hat{\mathbf{z}}_t$  se obține ca sumă ponderată cu coeficienții  $\alpha_{ii}$  a vectorilor de adnotare  $\mathbf{a}_i$ , multiplicată apoi cu scalarul  $\beta_t$  obținut prin aplicarea unei transformări liniare pe starea ascunsă anterioară a *decoder*-ului urmată de activare printr-o funcție de tip sigmoid(ă) – relația (14). Această codificare ponderată atențional  $\hat{\mathbf{z}}_t$  este concatenată cu *embedding*-ul cuvântului anterior (<start> la început)  $\mathbf{E}\mathbf{y}_{t-1}$ , formând intrarea de dimensiune  $2048 + 512 = 2560$  a LSTM, pe baza căreia este generată ca ieșire nouă stare ascunsă  $\mathbf{h}_t$  de dimensiune 512. Un strat liniar transformă această nouă stare ascunsă în scoruri de predicție asociate prin indice fiecărui cuvânt din vocabularul indexat disponibil, de dimensiune 9490, scoruri care sunt memorate. De asemenea sunt memorate și ponderile atenționale  $\alpha_{ii}$  la fiecare pas  $t$  spre a putea fi utilizate ulterior pentru vizualizarea pe imagini a regiunilor de focalizare dinamică a atenției la generarea fiecărui cuvânt  $\mathbf{y}_t$  din secvența de ieșire.

La antrenare, după fiecare epocă, poate fi salvat un fișier *checkpoint* cu valorile parametrilor modelului obținute până la momentul respectiv. Astfel, antrenarea poate fi reluată de la oricare

stadiu anterior al învățării cu parametrii modelului stocați într-un fișier *checkpoint* corespunzător. La sfârșitul fiecărei epoci de antrenare se execută o validare. Funcția obiectiv (*loss*) utilizată în cazul nostru de generare a unor secvențe de cuvinte este de tip *cross entropy*. Aceasta efectuează operațiile necesare pentru calcularea *softmax* și *log* pe baza scorurilor brute furnizate de stratul final al *decoder*-ului. La calcularea *loss*-ului, la fiecare pas de timp sunt ignorate regiunile completate cu  $\langle \text{pad} \rangle$ . În plus, așa cum am menționat mai sus, este utilizată o „regularizare dublu stocastică”, implicând și o a doua funcție *loss*, pentru modelul atențional. Astfel, se impune și ca suma tuturor ponderilor atenționale  $\alpha_{ii}$  ( $0 \leq \alpha_{ii} \leq 1$  și  $\sum_{i \leq L} \alpha_{ii} = 1$ ) pentru fiecare locație  $i$  să fie 1 pe totalitatea pașilor de timp  $T$  ( $\sum_{i \leq T} \alpha_{ii} \approx 1$ ), ceea ce impune modelului să se focalizeze pe fiecare pixel în cursul generării întregii secvențe descriptive de ieșire. Prin urmare se caută și minimizarea diferenței între 1 și suma ponderilor asociate pixelilor de-a lungul tuturor pașilor de timp – relația (15).

Pentru evaluarea performanței pe setul de validare și alegerea modelului (dar și pentru testare) se utilizează metrica automatizată BLEU (*BiLingual Evaluation Understudy*) 1 – 4. Aceasta a constituit multă vreme un standard de evaluare în raportările din literatura de specialitate în domeniul generării automate de fraze la traducerea automată dintr-un limbaj natural în altul, dar și de descrieri textuale asociate imaginilor. Scorul BLEU indică gradul de potrivire – considerând secvențe de 1 până la 4 cuvinte – cu un text de referință corespunzător creat de om și poate avea valori reale în intervalul  $[0, 1]$ , unde valoarea 1 indică potrivirea perfectă, iar valoarea 0 totală nepotrivire. Algoritmii BLEU tind prin construcție să acorde scorurile cele mai mari textelor generate cu lungimi apropiate de cele ale textelor de referință. În cazul nostru, pentru fiecare descriere generată, sunt utilizate toate cele  $N_c$  descrieri disponibile ca referință pentru respectiva imagine. După un anumit număr de epoci se constată apariția unei fracturi de corelare între continuarea încercărilor de minimizare a funcției *loss* de tip probabilitate logaritmică și evoluția scorului BLEU. Astfel, antrenarea trebuie oprită anticipat (*early stopping*) atunci când scorul BLEU începe să se degradeze chiar dacă *loss*-ul continuă să scadă. Această oprire anticipată a învățării pe baza scorului BLEU este singura strategie de regularizare utilizată în afară de *dropout*. A fost utilizată implementarea BLEU din modulul de resurse și instrumente Python pentru limbaj natural NLTK (*Natural Language Tool Kit*). Deoarece există destule critici aduse scorului BLEU, constatându-se că nu este corelat întotdeauna bine cu modul de gândire uman, [Xu et al., 2015] au utilizat de asemenea ca metrică alternativă scorul METEOR [Denkowski & Lavie, 2014], mult mai potrivit din acest punct de vedere, dar acesta nu a fost folosit în implementarea experimentată. Ca strategie de antrenare este recomandată abordarea învățării în etape. Astfel, se pornește mai întâi cu antrenarea *decoder*-ului fără ajustarea fină (*fine-tuning*) a *encoder*-ului, cu o dimensiune a loturilor (*batch size*) mai mare, pentru un număr de epoci  $M$ , observându-se evoluția scorului BLEU-4. Pentru optimizare se utilizează *optimizer*-ul Adam cu rată de învățare (*learning rate*) inițială de  $4e^{-4}$ . Se alege un *checkpoint* aferent unei epoci  $M_1$  anterioare  $M$ , pentru care scorul BLEU-4 a atins o valoare de circa 22,5-23,5. Pornind de la *checkpoint*-ul corespunzător acestei epoci  $M_1$ , se continuă antrenarea *encoder*-ului permițând și *fine-tuning*, alegând însă de această dată o dimensiune mai mică a loturilor, deoarece dimensiunea modelului include acum și gradientii *encoder*-ului. Se observă că la câteva epoci odată scorul BLEU-4 crește cu peste o unitate.

Atât la antrenare cât și la validare este aplicată tehnica *Teacher Forcing*, care constă în furnizarea ca intrare la fiecare nou pas de decodare a cuvântului real din descrierea curentă asociată, în locul cuvântului care tocmai a fost generat la pasul anterior. Utilizarea acestei tehnici este comună în timpul antrenării pentru accelerarea convergenței modelului. Totuși, acesta ar putea să devină astfel prea dependent de indicarea răspunsului corect și să ajungă să manifeste o oarecare instabilitate în practică. Ideal ar fi să se utilizeze *Teacher Forcing* la antrenare doar uneori, pe bază probabilistică (*scheduled sampling*). Pe de altă parte, în timpul validării ar trebui reproduse în cât mai mare măsură condițiile reale pentru inferență, deoarece cu *Teacher Forcing* și la validare, scorul BLEU calculat pentru descrierile astfel generate nu reflectă de fapt performanța reală a modelului. Practic, oprirea anticipată pe baza scorului BLEU ar trebui să fie utilizată doar dacă pentru validare este implementată inferența utilizând cuvântul anterior efectiv generat de model.

Mai trebuie menționat faptul că, pentru ajustarea fină a parametrilor modelului la *Transfer Learning*, este recomandat să se utilizeze o rată de învățare substanțial mai mică decât cea utilizată la antrenarea modelului original, pentru ca modificările parametrilor deja optimizați ai acestuia să nu fie prea grosiere și mai mult să strice. Optimizarea *encoder*-ului se face de asemenea cu Adam, dar cu o rată de învățare (*learning rate*) de  $1e^{-4}$ , ceea ce reprezintă o zecime din valoarea implicită pentru acest *optimizer*. Antrenarea unei epoci cu ajustarea fină a parametrilor/ponderilor arhitecturii (*fine-tuning*) poate dura în medie de circa 2,5-3 ori mai mult decât fără ajustare fină.

Prin vizualizarea efectului mapării ponderilor atenționale  $\alpha_{ti}$  pe diferite regiuni ale imaginilor de intrare la generarea fiecărui nou cuvânt în șirul descriptiv de ieșire la momentul  $t$ , se poate observa focalizarea dinamică selectivă a atenției și faptul că aceasta corespunde în mare măsura intuiției umane, adăugându-se astfel un nivel suplimentar de interpretabilitate a ieșirii modelului. Se poate observa faptul că mecanismul atențional „*soft*” permite focalizarea dinamică inclusiv pe regiuni semnificative care nu reprezintă obiecte.

## 4. Concluzii

Înțelegerea scenelor reprezintă încă un obiectiv ambițios în domeniul vederii artificiale. În acest context se încadrează și *Image Captioning* (IC) – generarea automată a unei descrieri textuale asociate unei imagini – cu identificarea obiectelor care apar în imagini și a atributelor lor, precum și a relațiilor și interacțiunilor dintre ele și exprimarea tuturor acestora în text în limbaj natural semnificativ și coerent.

Exemple de aplicații sunt: adnotarea automată cu texte descriptive a volumului tot mai mare de imagini postate pe rețelele de *social-media*; crearea de rezumate asociate imaginilor din diferite baze de date pentru organizare și/sau accesibilitate mai eficientă; asistarea persoanelor cu deficiențe vizuale majore prin descrierea textuală, combinată cu un sistem *text-to-speech* (TTS), a unor astfel de imagini, filme, sau chiar a mediului înconjurător în care se află sau se deplasează acestea; adnotarea cu descrieri ale celor observate în imagistica medicală de diferite tipuri pentru asistarea în punerea unor diagnostice și/sau urmărirea unor evoluții ș.a.m.d.

Odată cu evoluțiile recente spectaculoase din domeniul rețelelor neuronale adânci și *Deep Learning*, cu performanțe comparabile și chiar peste cele umane în domenii ca analiza de imagini și traducerea automată, au fost înregistrate progrese absolut remarcabile și în domeniul IC, calitatea textelor descriptive generate automat pentru imagini fiind îmbunătățită substanțial. În general, mai întâi se obțin reprezentări vectoriale esențializate ale caracteristicilor (*features*) din imagini folosindu-se rețele neuronale convoluționale, apoi aceste reprezentări sunt decodate în fraze în limbaj natural utilizându-se rețele neuronale recurente, în particular de tip LSTM.

Provocările actuale cele mai importante în domeniul IC se referă la generarea automată de descrieri textuale ale imaginilor cât mai apropiate de cele produse de subiecți umani, cât mai nuanțate și bogate în informații utile și care să se refere inclusiv la obiecte și situații nemaîntâlnite anterior la învățare. Există deja abordări și soluții cu rezultate promițătoare inclusiv pe aceste problematice. În lucrarea de față au fost prezentate sintetic principalele tehnici și metode bazate pe *Deep Learning* utilizate în ultimii ani în domeniul IC, cu detalierea și analizarea ca studiu de caz a uneia dintre cele mai performante dintre acestea, prima în care s-a utilizat (și) un mecanism atențional pe o arhitectură de tip *encoder-decoder* și pentru care s-a experimentat o implementare PyTorch.

## Mențiuni

Prezenta lucrare are la bază parte din activitățile și rezultatele etapei a II-a a proiectului PN 1937-0601 derulat la ICI București (2019) în cadrul Programului național nucleu „SMARTIC” finanțat de Ministerul Cercetării și Inovării.

**BIBLIOGRAFIE**

1. Anderson, P.; He, X.; Buehler, C.; Teney, D.; Johnson, M.; Gould, S. & Zhang, L. (2018). Bottom-up and top-down attention for image captioning and visual question answering. *CVPR2018*.
2. Bahdanau, D.; Cho, K. & Bengio, Y. (2014). Neural machine translation by jointly learning to align and translate. *arXiv:1409.0473*.
3. Cho, K.; van Merriënboer, B.; Gulcehre, C.; Bougares, F.; Schwenk, H. & Bengio, Y. (2014). Learning phrase representations using RNN encoder-decoder for statistical machine translation. In *EMNLP*, Oct. 2014.
4. Corbetta, M. & Shulman, G. L. (2002). Control of goal-directed and stimulus-driven attention in the brain. *Nature reviews neuroscience*, 3(3):201–215.
5. Denkowski, M. & Lavie, A. (2014). Meteor universal: Language specific translation evaluation for any target language. In *Proceedings of the EACL 2014 Workshop on Statistical Machine Translation*.
6. He, K.; Zhang, X.; Ren, S. & Sun, J. (2015). Deep Residual Learning for Image Recognition. Tech Report. *arXiv:1512.03385*.
7. He, S.; Tavakoli, H. R.; Borji, A. & Pugeault, N. (2019). A synchronized multi modal attention-caption dataset and analysis. *arXiv:1903.02499*.
8. Hochreiter, S. & Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, 9(8):1735–1780.
9. Hossain, M. D.; Soheli, F.; Shiratuddin, M. F. & Laga, H. (2019). A comprehensive survey of deep learning for image captioning. *ACM Computing Surveys (CSUR)*, 51(6):118.
10. Kingma, D. P. & Ba, J. (2014). Adam: A Method for Stochastic Optimization. *arXiv:1412.6980*.
11. Lu, J.; Xiong, C.; Parikh, D. & Socher, R. (2016). Knowing When to Look: Adaptive Attention via A Visual Sentinel for Image Captioning. *IEEE Conference on Computer Vision and Pattern Recognition – CVPR 2017*.
12. Mishra, A. & Liwicki, M. (2019). Using deep object features for image descriptions. *arXiv:1902.09969*.
13. Pascanu, R.; Gulcehre, C.; Cho, K. & Bengio, Y. (2014). How to construct deep recurrent neural networks. In *ICLR 2014*.
14. Ren, S.; He, K.; Girshick, R. & Sun, J. (2015). Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. *arXiv:1506.01497*.
15. Rensink, R. A. (2000). The dynamic representation of scenes. *Visual cognition*, 7(1-3):17–42.
16. Simonyan, K. & Zisserman, A. (2014). Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556.
17. Srivastava, N.; Hinton, G.; Krizhevsky, A.; Sutskever, I. & Salakhutdinov, R. (2014). Dropout: A simple way to prevent neural networks from overfitting. *JMLR 15*.
18. Staniute, R. & Sesok, D. (2019). A Systematic Literature Review on Image Captioning. *Applied Sciences*, 9(10):2024.
19. Sutskever, I.; Vinyals, O. & Le, Quoc V.V. (2014). Sequence to sequence learning with neural networks. In *NIPS 2014*:3104–3112.
20. Tieleman, T. & Hinton, G. (2012). Lecture 6.5 - RMSprop. *Technical report*.
21. Venugopalan, S.; Hendricks, L. A.; Rohrbach, M.; Mooney, R. J.; Darrell, T. & Saenko, K. (2017). Captioning Images with Diverse Objects. *Proc. of IEEE Conference on Computer Vision and Pattern Recognition – CVPR 2017*:1170-1178.



22. Vinodababu, S. (2018). PyTorch Tutorial to Image Captioning. <https://github.com/sgrvinod/a-PyTorch-Tutorial-to-Image-Captioning> (accesat 2019)
23. Vinyals, O.; Toshev, A.; Bengio, S. and Erhan, D. (2014). Show and tell: A neural image caption generator. *arXiv:1411.4555*.
24. Vrejoiu, M. H. (2019a). Rețele neuronale convoluționale, Big Data și Deep Learning în analiza automată de imagini. *Revista Română de Informatică și Automatică*, 29(1):91-114.
25. Vrejoiu, M. H. (2019b). Neural Networks and Deep Learning in Cyber Security. *Romanian Cyber Security Journal*, 1(1):69-86.
26. Xu, K.; Ba, J.; Kiros, R.; Cho, K.; Courville, A.; Salakhutdinov, R.; Zemel, R. & Bengio, Y. (2015). Show, Attend and Tell: Neural Image Caption Generation with Visual Attention. In Proceedings of the 32nd International Conference on Machine Learning, *Proc. of ML Research*, 37:2048-2057. (Preprint *arXiv:1502.03044*).
27. Wang, C.; Yang, H.; Bartz, C. & Meinel, C. (2016). Image captioning with deep bidirectional LSTMs. *Proc. of ACM on Multimedia Conference – ACM 2016*:988–997.



**Anca Mihaela HOTĂRAN** este cercetător științific gradul III în Departamentul „Sisteme inteligente distribuite intensive ca date” din ICI București. Domeniile și subiectele sale principale de expertiză și interes cuprind: reprezentarea cunoștințelor (ontologii, semantică), analiza avansată a datelor și inteligența artificială (sisteme expert, învățare automată).

**Anca Mihaela HOTĂRAN** is a 3<sup>rd</sup> grade scientific researcher in the Department of “Distributed Data Intensive Intelligent Systems” at ICI Bucharest. Her main areas and topics of expertise and interest include: Knowledge Representation (ontologies, semantics), advanced Data Analysis and Artificial Intelligence (Expert Systems, Machine Learning).



**Mihnea Horia VREJOIU** este cercetător științific gradul III în Departamentul „Sisteme inteligente distribuite intensive ca date” din ICI București. Domeniile și subiectele sale principale de expertiză și interes cuprind: vedere artificială (prelucrare și analiză de imagini, recunoașterea formelor, recunoașterea optică de caractere – OCR, recunoașterea numerelor de înmatriculare – LPR) și învățare automată (clasificatoare, memorii asociative).

**Mihnea Horia VREJOIU** is scientific researcher 3<sup>rd</sup> grade in the “Distributed Data Intensive Intelligent Systems” Department at ICI Bucharest. His main areas and topics of expertise and interest cover: Artificial Vision (Image Processing and Analysis, Pattern Recognition, Optical Character Recognition – OCR, License Plate Recognition – LPR) and Machine Learning (classifiers, associative memories).