

Inteligencia Artificial

Una Revolución

Pablo Zegers, Ph.D.
pablozegers@gmail.com

2019

Presentación

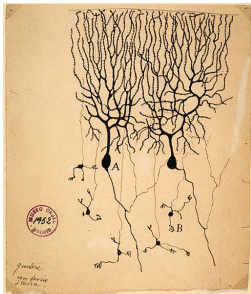
Científico y Emprendedor, Enfocado en Inteligencia Artificial



- No es una ciencia más, es una **meta ciencia**.
- **Inteligencia Artificial**
 - Aprendizaje de Máquinas
 - Aprendizaje Profundo (*Deep Learning*)
- **Robótica**, IA en un sistema físico capaz de interactuar con el mundo físico.

Inteligencia Artificial

Una Historia que Partió al Inicio del Siglo XX



Ramón y Cajal, *Revista Trimestral de Histología Normal y Patológica*, Laboratorio de Histología de la Facultad de Medicina de Barcelona, **1888**. Dibujo del cerebelo de un pichón, hecho por Santiago Ramón y Cajal (Wikipedia).

BULLETIN OF
MATHEMATICAL BIOPHYSICS
VOLUME 5, 1943

A LOGICAL CALCULUS OF THE IDEAS IMMANENT IN NERVOUS ACTIVITY

WARREN S. MCCULLOCK AND WALTER PITTS

FROM THE UNIVERSITY OF ILLINOIS, COLLEGE OF MEDICINE,
DEPARTMENT OF PSYCHIATRY AT THE ILLINOIS NEUROPSYCHIATRIC INSTITUTE,
AND THE UNIVERSITY OF CHICAGO

Because of the "all-or-none" character of nervous activity, neural events and the relations among them can be treated by means of propositional logic. It is found that the behavior of every net can be described in these terms, with the addition of some specialized logical means for nets containing circles; and that for any logical expression qualifying certain conditions, one can find a net behaving in the fashion it describes. It is shown that many particular circuits among possible neurophysiological structures are equivalent, in the sense that, for every, and behaving under one assumption, there exists another on which behaves under the other and gives the same results, although perhaps not in the same time. Various applications of the calculus are discussed.

McCullock and Pitts, *Journal of Mathematical Biophysics*, **1943**.

$$y = g \left(\sum_{i=1}^n w_i \cdot x_i - b \right)$$

$$g(s) = \frac{2}{1 + e^{-a \cdot s}}$$

Rosenblatt, *Principles of Neurodynamics: Perceptrons and the Theory of Brain Mechanisms*, Spartan Books, **1962**.

- A partir del 2011, sorpresivamente y en contra del sentido común, las **redes neuronales profundas resolvieron problemas difíciles con facilidad**.
- Schwartz-Ziv and Tishby, el 29 of April of 2017, **demuestran en “Opening the black box of Deep Neural Networks via Information”, ArXiv, que se cruzó un umbral matemático, hasta el momento desconocido, que permitió diseñar con facilidad IAs complejas.**

Inteligencia Artificial

Estado del Arte

arXiv:1611.08235v2 [cs.CV] 6 Dec 2016

Full-Resolution Residual Networks for Semantic Segmentation in Street Scenes

Tobias Pohlen, Alexander Hermans, Markus Matthias, Bastian Leibe

Visual Computing Institute

RWTH Aachen University

tobias.pohlen@vis-vc.informatik.rwth-aachen.de, {hermans, matthias, leibe}@informatik.rwth-aachen.de

Abstract

Semantic image representation is an essential component of modern autonomous driving systems, as an accurate understanding of the surrounding scene is crucial to navigation and action planning. Current state-of-the-art approaches to semantic image segmentation rely on pre-trained networks that were initially developed for classifying images as a whole. While these networks exhibit outstanding recognition performance (i.e., what is visible?), they lack localization accuracy (i.e., where precisely is something located?). Therefore, additional processing steps have to be performed in order to obtain pixel-accurate segmentation results at the full image resolution. To achieve this, we propose a novel ResNet-like architecture that exhibits strong localization and recognition power. We combine multi-scale context with pixel-level accuracy by using two processing streams: while our network does semantic segmentation at the full image resolution, enabling precise adherence to object boundaries. The other stream undergoes a sequence of pooling operations to obtain coarse features for recognition. The two streams are coupled at the full image resolution using residual units. Without additional pre-training and without pre-training, our approach achieves an intersection-over-union score of 77.3% on the Cityscapes dataset.

1. Introduction

Recent years have seen an increasing interest in self-driving cars and in driver assistance systems. A crucial aspect of autonomous driving is to acquire a comprehensive understanding of the surroundings in which a car is moving. Semantic image segmentation [1, 2, 3, 4, 5, 6], the task of assigning a set of predefined class labels to image pixels, is an important tool for modeling the complex relationships of the semantic content usually found in street scenes, such as cars, pedestrians, road, or sidewalks. In autonomous navigation it is used to estimate metrics, e.g., to open processing steps to discard image regions that are unlikely to contain objects of interest [7, 8] to improve object detection [9, 10, 11, 12].



(a) Pooling (b) Residual stream

(c) Pooling (d) Unpooling

(e) Residual stream (f) Pooling stream

(g) Unpooling (h) Residual stream

(i) Pooling (j) Residual stream

(k) Unpooling (l) Residual stream

(m) Pooling (n) Residual stream

(o) Unpooling (p) Residual stream

(q) Pooling (r) Residual stream

(s) Unpooling (t) Residual stream

(u) Pooling (v) Residual stream

(w) Unpooling (x) Residual stream

(y) Pooling (z) Residual stream

(aa) Unpooling (ab) Residual stream

(ac) Pooling (ad) Residual stream

(ae) Unpooling (af) Residual stream

(ag) Pooling (ah) Residual stream

(ai) Unpooling (aj) Residual stream

(ak) Pooling (al) Residual stream

(am) Unpooling (an) Residual stream

(ao) Pooling (ap) Residual stream

(aq) Unpooling (ar) Residual stream

(as) Pooling (at) Residual stream

(au) Unpooling (av) Residual stream

(aw) Pooling (ax) Residual stream

(ay) Unpooling (az) Residual stream

(ba) Pooling (bb) Residual stream

(bc) Unpooling (bd) Residual stream

(be) Pooling (bf) Residual stream

(bg) Unpooling (bh) Residual stream

(bi) Pooling (bj) Residual stream

(bk) Unpooling (bl) Residual stream

(bm) Pooling (bn) Residual stream

(bo) Unpooling (bp) Residual stream

(bq) Pooling (br) Residual stream

(bs) Unpooling (bt) Residual stream

(bu) Pooling (bv) Residual stream

(bw) Unpooling (bx) Residual stream

(by) Pooling (bz) Residual stream

(ca) Unpooling (cb) Residual stream

(cc) Pooling (cd) Residual stream

(ce) Unpooling (cf) Residual stream

(ca) Pooling (cb) Residual stream

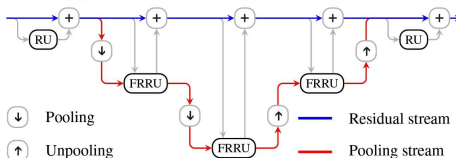


Figure 1. Example output and the abstract structure of our full-resolution residual network. The network has two processing streams. The residual stream (blue) stays at the full image resolution, the pooling stream (red) undergoes a sequence of pooling and unpooling operations. The two processing streams are coupled using full-resolution residual units (FRRUs).

2016-12-6

Deep Video Portraits

HYUNGWOO KIM, Max Planck Institute for Informatics, Germany
PABLO CARRERO, Technical Fellow
AYUSH Tewari and WEIPENG XU, Max Planck Institute for Informatics, Germany
JUDITH TREB and MATTHIAS NIESNER, Technical University of Munich, Germany
PATRICK PÉREZ, Technical Fellow
CHRISTIAN RICHARDT, University of Bath, United Kingdom
MICHAEL ZOLLHOFFER, Stanford University, United States of America
CHRISTIAN THEOBALT, Max Planck Institute for Informatics, Germany

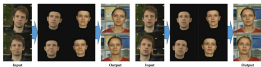


Fig. 1. Unlike current face reenactment approaches that only modify the expression of a target actor in a video, our novel deep video portrait approach enables full-frame-to-full-frame transfer of the target’s head pose, facial expression and eye motion with a high level of precision.

We present a novel approach that enables photo-realistic re-animation of portrait videos using only an input video. In contrast to existing approaches that require reconstruction of facial expressions only, we use the fact that humans do full-face movements and control our expressions on a pixel-by-pixel basis from camera space in a general video of target actor. The core of our approach is a generative neural network with several space-time autoencoders. The network takes as input spatial-temporal patches of a source face and, based on which it generates photo-realistic video frames for a given target actor. The realism in this re-animating is video transfer as learned by spatial-temporal training, such as a weak-to-strong transfer learning scheme that first trains the neural network—then taking full control of the

target. With the ability to freely combine source and target parameters, we are able to demonstrate a large variety of video transfer applications without explicitly modifying face, body or background. In addition, we use our model for full-face movements on cross-modal settings and transfer high fidelity video dubbing. To demonstrate the high quality of our output, we conduct an extensive series of experiments and evaluations. We also present a user study that shows that our video-clip are hard to detect.

CC BY 4.0 International license.

Additional Key Words and Phrases: Facial Reanimation, Video Portraits, Dubbing, Deep Learning, Conditional GAN, Clustering in Video Transitions

ACM Reference Format:

Hyungwoo Kim, Pablo Carrero, Ayush Tewari, Weipeng Xu, Judith Treb, Michael Zollhofer, Patrick Pérez, Christian Theobalt, and Matthias Niesner. 2018. Deep Video Portraits. In Proceedings of ACM SIGGRAPH 2018. ACM, New York, NY, USA, 14 pages. <https://doi.org/10.1145/3201251>

1 INTRODUCTION

Digitizing and editing video portraits, i.e. videos aimed to show a person’s head and upper body, is an important problem in computer graphics, with applications in video dubbing and movie post-production, visual effects, virtual dubbing, virtual reality, and tele-presence among others. In this paper, we address the problem of re-animating a photo-realistic video portrait of a target actor that captures the actions of a source actor whose expression and target actor’s face are full-control of the target actor.

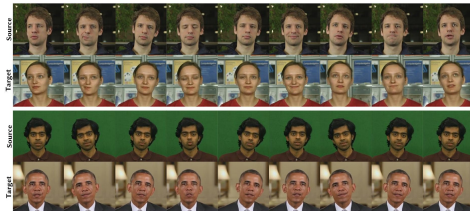


Fig. 5. Qualitative results of full-head reenactment: our approach enables full-frame target video portrait synthesis under full 3D head pose control. The output video portraits are photo-realistic and hard to distinguish from real videos. Note that even the shadow in the background of the second row moves consistently with the modified foreground head motion. In the sequence at the top, we only transfer the translation in the camera plane, while we transfer the full 3D translation for the sequence at the bottom. For full sequences, please refer to our video. Obama video courtesy of the White House (public domain).

ACM SIGGRAPH 2018, July 29–30, 2018, Houston, Texas, USA

2018-5-29

Learning Dexterous In-Hand Manipulation



Figure 1: A blue fingered humanoid hand trained with reinforcement learning manipulating a block from an initial configuration to a goal configuration using vision for sensing.

Abstract

We use reinforcement learning (RL) to learn dexterous in-hand manipulation policies which can perform vision-based object manipulation on a physical Shadow Hand. The training is performed in a simulated environment to learn a combination many of the physical properties of the system like friction coefficients and an object's appearance. Our policies transfer to the physical robot despite being trained entirely in simulation. Our method does not rely on any human demonstrations, but more reliance toward human manipulation energy naturally including finger posing, wrist finger coordination, and the controlled use of gravity. Our results were obtained using the same distributed RL system that was used to train OpenAI Five [15]. We also include a video of our results. Images: /papers/167/abstract.html

1 Introduction

While dexterous manipulation of objects is an fundamental everyday task for humans, it is still challenging for autonomous robots. Modern-day robots are typically designed for specific tasks in controlled settings and are largely unable to adjust, adapt, and effectively. In contrast, people are able to perform a wide range of dexterous manipulation tasks in a diverse set of environments, making the human hand a generalist source of inspiration for robotic in-hand manipulation. The Shadow Hand (SH) is an example of a robotic hand designed for human-level dexterity. It has five fingers with a total of 38 degrees of freedom. The hand has been commercially available since 2015 (<http://www.shadow-robot.com>) and is open to OpenAI in alphabetical order.

Shao-An Huang, Benoit Bonet, Minh-Chien Chen, Rishi Desai, Benoit Drouot, Shikhar Ghosh, Armin Greuter, Markéta Popović, Shuang Tang, Shi-Kei Tang, Antonio Torralba, Tomer Shalev-Shanzer, Pieter Abbeel, Lilian Weng, Wojciech Zarembki

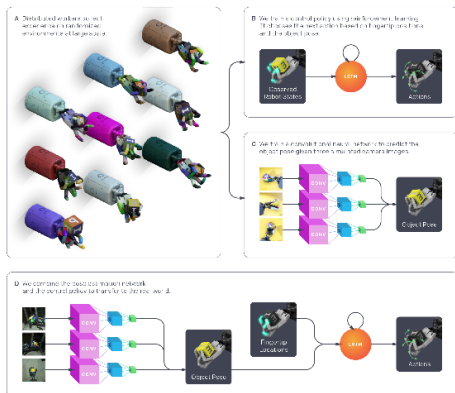


Figure 2: System Overview. (a) We use a large distribution of simulations with randomized parameters and appearances to collect data for both the control policy and vision-based pose estimator. (b) The control policy receives observed robot states and rewards from the distributed simulations and learns to map observations to actions using a recurrent neural network and reinforcement learning. (c) The vision based pose estimator renders scenes collected from the distributed simulations and learns to predict the pose of the object using a convolutional neural network (CNN), trained separately from the control policy. (d) To transfer to the real world, we predict the object pose from 3 real camera feeds with the CNN, measure the robot fingertip locations using a 3D motion capture system, and give both of these to the control policy to produce an action for the robot.

2018-7-30

Linguistic Regularities in Continuous Space Word Representations

Tomas Mikolov¹, Wen-tau Yih, Geoffrey Zweig
Microsoft Research
Redmond, WA 98052

Abstract

Continuous space language models have recently demonstrated outstanding results across a variety of tasks. In this paper, we examine the vector space word representations that are implicitly learned by the unsupervised weights. We find that these representations are surprisingly good at capturing syntactic and semantic regularities in language, and that such relationships are observed by a relation-specific vector offset. The above vector-oriented reasoning based on the offsets between words. For example, the multi-formal relationship is automatically learned, and with the learned vector representations, “King - Man + Queen” results in a vector very close to “Queen”. We demonstrate that the word co-occurrence system replicates the scores of syntactic analogy questions (provided with this paper), and is able to correctly answer almost 80% of the questions. We demonstrate that the word vectors capture semantic regularities by using the vector offset method to answer Stanford 2012 task 2 questions. Remarkably, the method outperforms the best previous system.

1 Introduction

A defining feature of neural network language models is their representation of words in high-dimensional vector spaces. In these models (Bengio et al., 2003; Mikolov, 2005; Mikolov et al., 2010), words are converted via a learned lookup table into real-valued vectors which are used as the

inputs to a neural network. As pointed out by the original proposers, one of the main advantages of these models is that the distributed representation achieves a level of generalization that is not possible with classical n-gram language models: whereas a n-gram model needs to learn a discrete vector that has an inherent relationship to one another, a continuous space model needs to learn a word vector where similar words are likely to have similar vectors. Thus, when the model parameters are adjusted in response to a particular word or word sequence, the improvements will carry over to occurrences of similar words and sequences.

By training a neural network language model, one obtains not just the model itself, but also the learned word representations, which may be used for other, potentially unrelated, tasks. This has been used to good effect, for example in Schuster and Weston, 2008; Tenen et al., 2010 where learned word representations are used with specialized classifiers to improve performance in many NLP tasks.

In this work, we find that the learned word representations in fact capture meaningful syntactic and semantic regularities in a very simple way. Specifically, the regularities are observed as constant vector offsets between pairs of words sharing a particular relationship. For example, if we denote the vector for word i as v_i , and focus on the straightforward relation, we observe that $v_{\text{Queen}} - v_{\text{Man}} = v_{\text{King}} - v_{\text{Man}} = v_{\text{Queen}} - v_{\text{King}}$ and so on. Perhaps more surprisingly, we find that this also the case for a variety of semantic relations, as measured by the Stanford 2012 task of measuring relation similarity.

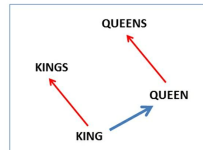
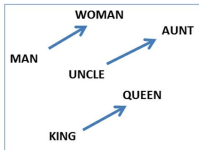


Figure 2: Left panel shows vector offsets for three word pairs illustrating the gender relation. Right panel shows a different projection, and the singular/plural relation for two words. In high-dimensional space, multiple relations can be embedded for a single word.

2013-6-10

Inteligencia Artificial

Estado del Arte

Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation

Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi, pingduo_jiang, satori, zhi-feng, qv, mnorouzi@google.com

Wolfgang Macherey, Miroslav Miklavic, Yanzhao Chen, Qin Ge, Klaus Moravets, Jeff Klingner, Apurva Shah, Mikko Lohman, Xinghui Liu, Lukasz Kaiser, Stephan Gouws, Yoshitaka Kato, Tamas Kocz, Hideo Koizumi, Keith Stevens, George Diner, Nikolaus Pelekh, Wei Wang, Cliff Young, Jason Riesa, Jason Riesa, Alex Buzhuk, Ovid Vinyals, Georgios Corrado, Marc'Antonio Frete, Jeffrey Dean

Abstract

Neural Machine Translation (NMT) is an end-to-end learning approach for automatic translation, with the potential to overcome many of the weaknesses of conventional phrase-based translation systems. Unfortunately, NMT systems are known to be computationally expensive both in training and in translation inference—sometimes prohibitively so in the case of very large data sets and large models. Several authors have also argued that NMT systems lack robustness, particularly when input sentences contain rare words. These issues have hindered NMT's use in practical deployments and services, where both accuracy and speed are essential. In this work, we present Google's Neural Machine Translation system, which attempts to address many of these issues. Our model consists of a deep LSTM network with a residual unit's flexible layers using residual connections as well as attention connections from the decoder network to the encoder. To improve translation and inference efficiency, we use an attention mechanism to compute the hidden layer of the decoder to use for each step of the encoder. In addition to the final translation step, we employ a low-precision softmax during inference computation. To improve handling of rare words, we divide words into a hierarchy of difficulty and model each word's "translation" for both input and output. This method provides a good balance between the flexibility of "abstract" identical words and the efficiency of "word" identical words, thereby enabling translation of rare words and ultimately improving the overall accuracy of the system. Our beam search decoder employs a length-normalization procedure and uses a coverage penalty, which encourages generation of shorter sentences that is most likely to cover all the words in the source sentence. We directly optimize the translation BLEU scores, we consider replacing the model by using reinforcement learning, but we found that the improvement in the BLEU scores did not reflect the human evaluation. On the WMT14 English-to-French and English-to-German benchmarks, Google's NMT system consistently results in state-of-the-art. Using a human-like evaluation on a set of translated English sentences, we observe translation errors for an average of 80% compared to Google's phrase-based production system.

1 Introduction

Neural Machine Translation (NMT) [1, 2] has recently been introduced as a promising approach with the potential of addressing many shortcomings of traditional machine translation systems. The strength of NMT lies in its ability to learn, directly, an end-to-end function, the mapping from input text to associated output text. Its architecture typically consists of two separate neural networks (NNs), one to generate the input text sequence and one to generate the associated output. NMT is often accompanied by an encoder-decoder architecture [2] which helps it cope effectively with long input sequences.

An advantage of Neural Machine Translation is that it achieves more flexible designs than traditional phrase-based machine translation [3]. In practice, however, NMT systems tend to be more computationally expensive than traditional translation systems, especially when training on very large-scale datasets or when serving the very best quality available translation services. These inherent weaknesses of Neural Machine Translation are

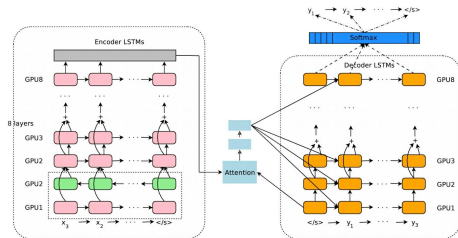


Figure 1: The model architecture of GNMT, Google's Neural Machine Translation system. On the left is the encoder network, on the right is the decoder network, in the middle is the attention module. The bottom encoder layer is bi-directional: the pink nodes gather information from left to right while the green nodes gather information from right to left. The other layers of the encoder are uni-directional. Residual connections start from the layer third from the bottom in the encoder and decoder. The model is partitioned into multiple GPUs to speed up training. In our setup, we have 8 encoder LSTM layers (1 bi-directional layer and 7 uni-directional layers), and 8 decoder layers. With this setting, one model replica is partitioned 8-ways and is placed on 8 different GPUs typically belonging to one host machine. During training, the bottom bi-directional encoder layers compute in parallel first. Once both finish, the uni-directional encoder layers can start computing, each on a separate GPU. To retain as much parallelism as possible during running the decoder layers, we use the bottom decoder layer output only for obtaining recurrent attention context, which is sent directly to all the remaining decoder layers. The softmax layer is also partitioned and placed on multiple GPUs. Depending on the output vocabulary size we either have them run on the same GPUs as the encoder and decoder networks, or have them run on a separate set of dedicated GPUs.

2016-10-8

A Deep Reinforced Model for Abstractive Summarization

Román Pardo, Caixiang Xiong and Richard Socher
(rpadilla, cxiong, rsocher}@jaxlabforce.com

Abstract

Attentional, RNN-based encoder-decoder models for abstractive summarization have achieved good performance on short input and output sequences. However, for longer documents and summaries, these models often include repetitive and incoherent phrases. We introduce a neural network model with reinforcement and a new training method. This method combines standard supervised word prediction and reinforcement learning (RL). Models trained only with the former often exhibit “stoppage bias” – they assume good truth is provided at each step during training. However, when standard word prediction is combined with the global reward provided by testing of RL, the resulting summaries become more readable. We evaluate this model on the CNN/Daily Mail and New York Times datasets. Our model obtains a 41.6 BLEU-4 score on the CNN/Daily Mail dataset, a 5.7 absolute points improvement over previous state-of-the-art models. It also performs well as the first abstractive model on the New York Times corpus. Human evaluation also shows that our model produces higher quality summaries.

1 Introduction

Text summarization is the process of automatically generating natural language summaries from an input document while retaining the important parts.

By condensing large quantities of information into short, informative summaries, summarization can aid many downstream applications such as

creating news digests, tweets, and report generators.

There are two prominent types of summarization algorithms. First, extractive summarization extracts facts summaries by copying parts of the input (Nee et al., 2002; Don et al., 2003; Nalapat et al., 2007). Second, abstractive summarization systems generate new phrases, possibly rephrasing or using words that were not in the original text (Chopra et al., 2016; Nalapat et al., 2016; Zeng et al., 2016).

Recently, neural network models (Nalapat et al., 2016; Zeng et al., 2016), based on the attentional encoder-decoder model for machine translation (Bahdanau et al., 2014), were able to generate abstractive summaries with high BLEU-4 scores. However, these systems have typically focused on summarizing short input sequences (one or two sentences) to generate even shorter summaries. For example, the summaries on the DUC-2000 dataset generated by the state-of-the-art system by Zeng et al. (2016) are limited to 75 characters.

Nalapat et al. (2016) also applied their abstractive summarization model to the CNN/Daily Mail dataset (Hermann et al., 2015), which contains input sequences of up to 500 tokens and multi-sentence summaries of up to 100 tokens. The analysis by Nalapat et al. (2016) illustrates a key problem with attentional encoder-decoder models: they often generate redundant summaries consisting of repeated phrases.

We present a new abstractive summarization model that achieves state-of-the-art results on the CNN/Daily Mail and similarly good results on the New York Times dataset (NVT) (Stanbury, 2006). To our knowledge, this is the first model for abstractive summarization on the NVT dataset. We introduce a key attention mechanism and a new training objective to address the repeating phrase

Source document

Jenson Button was denied his 100th race for McLaren after an ERS prevented him from making it to the start-line. It capped a miserable weekend for the Briton; his time in Bahrain plagued by reliability issues. Button spent much of the race on Twitter delivering his verdict as the action unfolded. 'Kimi is the man to watch', and 'loving the sparks', were among his pearls of wisdom, but the tweet which courted the most attention was a rather mischievous one: 'Ooh is Lewis backing his team mate into Vettel?' he quizzed after Rosberg accused Hamilton of pulling off such a manoeuvre in China. Jenson Button waves to the crowd ahead of the Bahrain Grand Prix which he failed to start Perhaps a career in the media beckons Lewis Hamilton has out-qualified and finished ahead of Nico Rosberg at every race this season. Indeed Rosberg has now beaten his Mercedes team-mate only once in the 11 races since the pair infamously collided in Belgium last year. Hamilton secured the 36th win of his career in Bahrain and his 21st from pole position. Only Michael Schumacher (40), Ayrton Senna (29) and Sebastian Vettel (27) have more. He also became only the sixth F1 driver to lead 2,000 laps. Nico Rosberg has been left in the shade by Lewis Hamilton who celebrates winning his third race of the year Kimi Raikkonen secured a record seventh podium finish in Bahrain following his superb late salvo, although the Ferrari driver has never won in the Gulf Kingdom. It was the Finn's first trip to the rostrum since the 2013 Korean Grand Prix, but his triumph brought a typically deadpan response: 'You're never happy when you finish second... I'm a bit pleased to get a result.' Sparks fly off the back of Kimi Raikkonen's Ferrari en route to finishing second in Bahrain Bernie Ecclestone was in the Bahrain paddock this weekend. He denied trying to engineer a deal for Hamilton, out of contract at the end of the season, to join Ferrari despite earlier insisting that such a move would be "great" for the sport. The 84-year-old also confirmed that F1 would be in Azerbaijan for the first time next year, even with concerns surrounding the country's human rights record. 'I think everybody seems to be happy,' Ecclestone said. 'There doesn't seem to be any big problem there. There's no question of it not being on the calendar. It's going to be another good race. Formula One supremo Bernie Ecclestone speaks to Nico Rosberg ahead of the Bahrain Grand Prix

Ground truth summary

Button denied 100th race start for McLaren after ERS failure. Button then spent much of the Bahrain Grand Prix on Twitter delivering his verdict on the action as it unfolded. Lewis Hamilton has out-qualified and finished ahead of Mercedes team-mate Nico Rosberg at every race this season. Bernie Ecclestone confirms F1 will make its bow in Azerbaijan next season.

ML₁ with intra-attention (ROUGE-1 41.58)

Button was denied his 100th race for McLaren. ERS prevented him from making it to the start-line. The Briton. He quizzed after Nico Rosberg accused Lewis Hamilton of pulling off such a manoeuvre in China. Button has been in Azerbaijan for the first time since 2013.

RL₁ with intra-attention (ROUGE-1 50.00)

Button was denied his 100th race for McLaren after an ERS prevented him from making it to the start-line. It capped a miserable weekend for the Briton. Button has out-qualified. Finished ahead of Nico Rosberg at Bahrain. Lewis Hamilton has. In 11 races. The race. To lead 2,000 laps. In. . . And. .

ML+RL₁ with intra-attention (ROUGE-1 44.00)

Button was denied his 100th race for McLaren. The ERS prevented him from making it to the start-line. Button has his team mate in the 11 races in Bahrain. He quizzed after Nico Rosberg accused Lewis Hamilton of pulling off such a manoeuvre in China.

Table 3: Example from the CNN/Daily Mail test dataset showing the outputs of our three best models after de-tokenization, re-capitalization, replacing anonymized entities, and replacing numbers. The ROUGE score corresponds to the specific example.

2017-5-19

Inteligencia Artificial

Estado del Arte

arXiv:1412.2306v2 [cs.CV] 14 Apr 2015

Deep Visual-Semantic Alignments for Generating Image Descriptions

Andrej Karpathy Li Fei-Fei
Department of Computer Science, Stanford University
{karpathy, feifei}@cs.stanford.edu

Abstract

We present a model that generates natural language descriptions of images and their regions. Our approach leverages distant supervision and their anchor descriptions to learn about the inter-modal correspondences between language and visual data. Our alignment model is based on a novel combination of Convolutional Neural Networks over image regions, bidirectional Recurrent Neural Networks over sentences, and a structured objective that aligns the two modalities through a multimodal embedding. We then describe a Multiresolution Neural Network architecture that uses the regional alignments to learn to generate natural descriptions of image regions. We demonstrate that our alignment model produces state-of-the-art results in several experiments on Flickr3K, Flickr30K and MSCOCO datasets. We then show that the generated descriptions significantly outperform textual baselines on both full images and on a new dataset of region-level annotations.



Figure 1. Motivation/Concept Figure: The model treats language as a rich label space and generates descriptions of image regions.

generating these descriptions of image regions (Fig. 1). The primary challenge towards this goal is in the design of a model that is rich enough to simultaneously reason about contents of images and their representation in the domain of natural language. Ideally, the model should be free of assumptions about specific hand-crafted templates, rules or capacities and instead rely on learning from the training data. The second, practical challenge is that datasets of image captions are available in large quantities on the Internet [1, 76, 77], but these descriptions analyze variations of several indices whose locations in the images are unknown. One can imagine that we can leverage these large image-sentence datasets by treating the sentences as weak labels, to which candidate regions of words correspond to some particular, but unknown location in the image. Our approach is to infer these alignments and use them to learn a generative model of descriptions. Concretely, our contributions are twofold:

• We develop a deep neural network model that infers the latent alignment between segments of sentences and the regions of the image that they describe.

1. Introduction

A quick glance at an image is sufficient for a human to point out and describe an immense amount of details about the visual scene [10]. However, the remarkable ability has proven to be a herculean task for our visual recognition models. The majority of previous work in visual recognition has focused on labeling images with a fixed set of visual concepts and prior progress has been achieved in these endeavors [4, 11]. However, while closed vocabularies of visual concepts contain a convenient modeling assumption, they are rarely intuitive when compared to the common sense of multi-observations that a human can compute.

Some pioneering approaches that address the challenge of generating image descriptions have been developed [70, 13]. However, these models rely only on learning natural visual concepts and sentence templates, which imposes limits on their variety. Moreover, the focus of these works has been on reducing complex visual scenes into a single sentence, which we consider to be an unnecessary restriction.

In this work, we strive to take a step towards the goal of

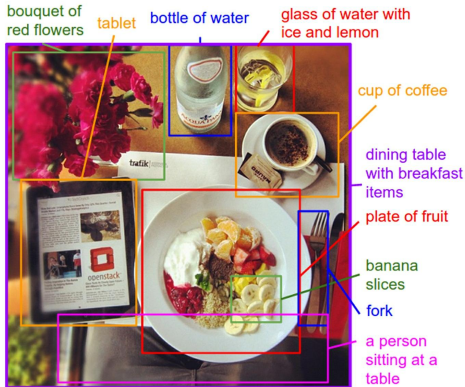


Figure 1. Motivation/Concept Figure: Our model treats language as a rich label space and generates descriptions of image regions.

StackGAN: Text to Photo-realistic Image Synthesis with Stacked Generative Adversarial Networks

Han Zhang¹, Tao Xu¹, Hongsheng Li², Shouling Zhang¹, Xiaohu Huang¹, Xiangqiang Wang¹, Dimitris Metaxas³

¹Department of Computer Science, Rutgers University

²Department of Computer Science and Engineering, Lehigh University

³Department of Electronic Engineering, The Chinese University of Hong Kong

⁴Department of Computer Science, University of North Carolina at Charlotte

Abstract

Synthesizing photo-realistic images from text descriptions is a challenging problem in computer vision and has many practical applications. Samples generated by existing text-to-image approaches can roughly reflect the meaning of the given descriptions, but they fail to contain accurate details and vivid object parts. In this paper, we propose stacked Generative Adversarial Networks (StackGAN) to generate photo-realistic images conditioned on text descriptions. The Stage-I GAN sketches the primitive shape and basic colors of the objects based on the given text description, yielding Stage-I low resolution images. The Stage-II GAN takes Stage-I results and text descriptions as inputs, and generates high resolution images with photo-realistic details. The Stage-II GAN is able to accept arbitrary and self-modifying details with the refinement process. Samples generated by StackGAN are more plausible than those generated by existing approaches. Apparently, our StackGAN is the first time generating realistic 256 × 256 images conditioned on only text descriptions, while state-of-the-art methods can generate at most 128 × 128 images. It demonstrates the effectiveness of the proposed StackGAN. Concrete experiments are conducted on CUB and Oxford-102 datasets, which contain enough object-agnostic variations and are widely used for text-to-image generation analysis.

1. Introduction

Generating photo-realistic images from text descriptions is a challenging problem. It has tremendous applications including photo editing and computer-aided image search.

¹hanzhang@rutgers.edu



Figure 1. Photo-realistic images generated by our StackGAN from unseen text descriptions. Descriptions for birds and flowers are from CUB [32] and Oxford-102 [18] datasets, respectively. (a) Given text descriptions, Stage-I of StackGAN sketches rough shapes and basic colors of objects, yielding low resolution images. (b) Stage-II of StackGAN takes Stage-I results and text descriptions as inputs, and generates high resolution images with photo-realistic details.

fully automatic synthesis systems are available. However, even the most advanced methods failed in generating high resolution images with photo-realistic details using text descriptions. The main challenge of this problem is that the space of plausible images given text description is multi-modal. There are a larger number of images that correspond to the given text descriptions.

Recently, Generative Adversarial Networks (GAN) [1, 14] have shown promising results in modeling complex multimodal data and generating multi-modal images. Reed et al. [15] demonstrated that GAN can effectively gener-

This bird has a yellow belly and tarsus, grey back, wings, and brown throat, nape with a black face

This bird is white with some black on its head and wings, and has a long orange beak

This flower has overlapping pink pointed petals surrounding a ring of short yellow filaments

(a) Stage-I images

(b) Stage-II images

Figure 1. Photo-realistic images generated by our StackGAN from unseen text descriptions. Descriptions for birds and flowers are from CUB [32] and Oxford-102 [18] datasets, respectively. (a) Given text descriptions, Stage-I of StackGAN sketches rough shapes and basic colors of objects, yielding low resolution images. (b) Stage-II of StackGAN takes Stage-I results and text descriptions as inputs, and generates high resolution images with photo-realistic details.

ARTICLE

Mastering the game of Go with deep neural networks and tree search

David Silver¹, Aja Huang¹, Clark B. White¹, Arthur Guez¹, Laurent Elie¹, George Van den Driessche¹, Marc-Aurèle Blais¹, Matthew H. Hoffman¹, John Borrajo¹, Sander Dieleman¹, Thore Graepel¹, Daniël J. Bernstein¹, Alexei A. Efros², Volodymyr Mnih¹, Demis Hassabis¹, David Silver¹, Markus Luck¹, Koray Kavukcuoglu¹, Thore Graepel¹ & Demis Hassabis¹

The game of Go has long been viewed as the most challenging of board games for artificial intelligence owing to its enormous search space and the difficulty of evaluating board positions and moves. Here we introduce a new approach to computer Go that uses a neural network to evaluate board positions and a policy network to select moves. These deep neural networks are trained by a novel combination of program search by learning from human expert games, and reinforcement learning from self-play. Without any knowledge of rules, Go board layout or rules of play, our program defeated the best human Go player in the world at a 9.5% margin of 100 games. This result was achieved by a program that played 100 million games of self-play, and defeated the human champion 100 times out of 100. This is the first time that a computer program has defeated a human professional player in the full-sized game of Go, a feat previously thought to be at least a decade away.

Artificial intelligence has made significant progress in a number of domains, such as image classification, speech recognition, and machine translation. However, it has struggled to master the game of Go, a board game that is widely regarded as the most challenging of board games for artificial intelligence. This is because of the enormous search space and the difficulty of evaluating board positions and moves. In this paper, we introduce a new approach to computer Go that uses a neural network to evaluate board positions and a policy network to select moves. These deep neural networks are trained by a novel combination of program search by learning from human expert games, and reinforcement learning from self-play. Without any knowledge of rules, Go board layout or rules of play, our program defeated the best human Go player in the world at a 9.5% margin of 100 games. This result was achieved by a program that played 100 million games of self-play, and defeated the human champion 100 times out of 100. This is the first time that a computer program has defeated a human professional player in the full-sized game of Go, a feat previously thought to be at least a decade away.

Artificial intelligence has made significant progress in a number of domains, such as image classification, speech recognition, and machine translation. However, it has struggled to master the game of Go, a board game that is widely regarded as the most challenging of board games for artificial intelligence. This is because of the enormous search space and the difficulty of evaluating board positions and moves. In this paper, we introduce a new approach to computer Go that uses a neural network to evaluate board positions and a policy network to select moves. These deep neural networks are trained by a novel combination of program search by learning from human expert games, and reinforcement learning from self-play. Without any knowledge of rules, Go board layout or rules of play, our program defeated the best human Go player in the world at a 9.5% margin of 100 games. This result was achieved by a program that played 100 million games of self-play, and defeated the human champion 100 times out of 100. This is the first time that a computer program has defeated a human professional player in the full-sized game of Go, a feat previously thought to be at least a decade away.

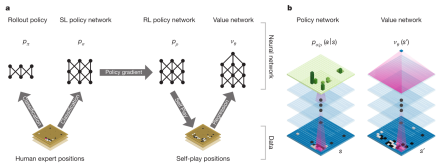


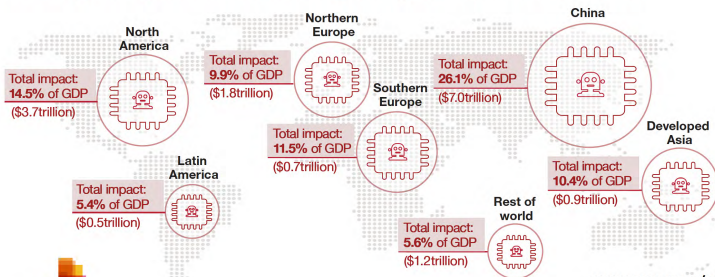
Figure 1 | Neural network training pipeline and architecture. a, A fast rollout policy p_1 and supervised learning (SL) policy network p_2 are trained to predict human expert moves in a data set of positions. A reinforcement learning (RL) policy network p_3 is initialized to the SL policy network, and is then improved by policy gradient learning to maximize the outcome (that is, winning more games) against previous versions of the policy network. A new data set is generated by playing games of self-play with the RL policy network. Finally, a value network v_1 is trained by regression to predict the expected outcome (that is, whether the current player wins) in positions from the self-play data set. b, Schematic representation of the neural network architecture used in AlphaGo. The policy network takes a representation of the board position s as its input, passes it through many convolutional layers with parameters θ (SL policy network) or μ (RL policy network), and outputs a probability distribution $p_\theta(a|s)$ or $p_\mu(a|s)$ over legal moves a , represented by a probability map over the board. The value network similarly uses many convolutional layers with parameters θ , but outputs a scalar value $v_\theta(s')$ that predicts the expected outcome in position s' .

2016-01-18

Inteligencia Artificial

Creación de US \$15,7 Millones de Millones de Valor Económico

Sizing the prize – Which regions gain the most from AI?



www.pwc.com/ai
#AIrevolution

© 2017 PricewaterhouseCoopers LLP. All rights reserved.

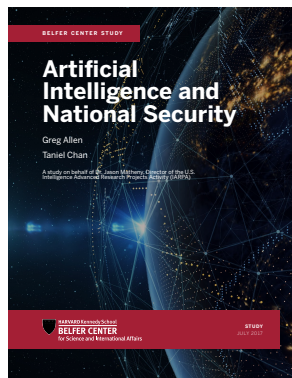
Inteligencia Artificial

La Fiebre de la Inteligencia Artificial



Personas caminando por el Sendero Chilkoot durante la Fiebre del Oro en Yukon, Alaska, a fines del siglo XIX (Wikipedia).

- Incide en **todas las actividades humanas**, no es sobre un producto o servicio.
- Afecta a **donde sea que haya personas**, no está localizada geográficamente.
- Produce una **reacción de empresas y gobiernos** que tiene como fin asegurar posiciones en este nuevo escenario.



● Executive summary

- Researchers in the field of Artificial Intelligence (AI) have demonstrated significant technical progress over the past five years, **much faster** than was previously anticipated.
- Most AI research advances are occurring in the **private sector and academia**.
- **Existing capabilities in AI have significant potential for national security.**
- Future progress in AI has the potential to be a transformative national security technology, **on a par with nuclear weapons, aircraft, computers, and biotech.**
- Advances in AI will affect national security by driving change in three areas: **military superiority, information superiority, and economic superiority.**
- We analyzed four prior cases of transformative military technologies—nuclear, aerospace, cyber, and biotech—and generated **"lessons learned" for AI.**
- Taking a "whole of government" frame, **we provide three goals for U.S. national security policy toward AI technology and provide 11 recommendations.**

Reacciones de los Gobiernos

Fuerzas Armadas



Reacciones de los Gobiernos

Capacidad de Manipulación

Psychological targeting as an effective approach to digital mass persuasion

L. C. Messa¹, M. Rounis^{1,2*}, G. Nava¹, and J. S. DellaVigna^{1,3}

¹Yale University School of Management, New Haven, CT 06520; ²Yale School of Journalism, New Haven, CT 06520; ³Yale School of Public Health, New Haven, CT 06520; *Correspondence: m.rounis@yale.edu

October 28, 2019; revised November 11, 2019; accepted November 11, 2019

People are exposed to persuasive communication across many contexts, including news, advertising, and political persuasion. It is important to understand how these communications affect behavior, and to identify ways to improve them. We study the effects of psychological targeting on persuasion in an online setting. We find that psychological targeting is effective in increasing persuasion across a wide range of products and services, and that this effect is particularly strong for individuals with high levels of extraversion. Psychological targeting is also effective in increasing persuasion across a wide range of products and services, and that this effect is particularly strong for individuals with high levels of extraversion. Psychological targeting is also effective in increasing persuasion across a wide range of products and services, and that this effect is particularly strong for individuals with high levels of extraversion.

Significance
Building on recent advances in the treatment of people based on their digital footprint, this paper demonstrates the effectiveness of psychological targeting—i.e., the use of persuasion appeals to the psychological characteristics of large groups of individuals with the goal of influencing their behavior. On the one hand, this form of psychological targeting is highly effective in increasing persuasion across a wide range of products and services. On the other hand, it could be used to target individuals with high levels of extraversion and low levels of openness to experience. This paper also discusses the implications of these findings for policy and practice.

10047819 | PNAS | November 28, 2019 | vol. 116 | no. 48 | 11048

A High Extraversion



Dance like no one's watching
(but they totally are)

Low Extraversion



Beauty doesn't have to shout

B High Openness



Aristoteles? The Seychelles? Unleash your creativity and challenge your imagination with an unlimited number of crossword puzzles!

Low Openness

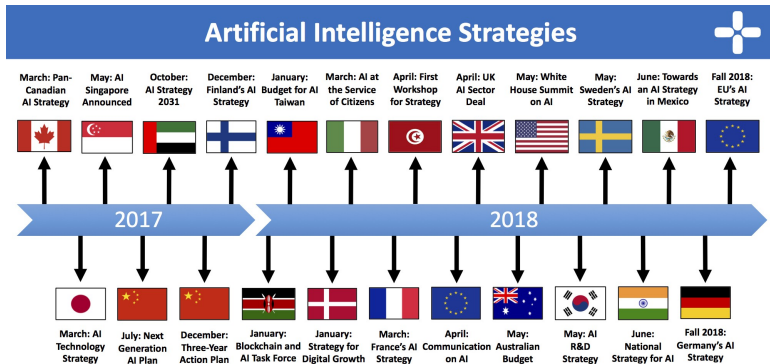


Settle in with an all-time favorite! The crossword puzzle that has challenged players for generations.

- **"Some countries are already moving in this direction. China has begun to construct a digital authoritarian state by using surveillance and machine learning tools to control restive populations, and by creating what it calls a "social credit system." Several like-minded countries have begun to buy or emulate Chinese systems. Just as competition between liberal democratic, fascist, and communist social systems defined much of the twentieth century, so the struggle between liberal democracy and digital authoritarianism is set to define the twenty-first."**
(How Artificial Intelligence Will Reshape the Global Order, Foreign Affairs, 10 de julio de 2018)

Reacciones de los Gobiernos

Planes Gubernamentales



2018-07-13 | Politics + AI | Tim Dutton

Reacciones de los Gobiernos

Objetivos Generales

- Aprovechar la **oportunidad**.
- Crear **capacidades humanas**.
- Ayudar a la **transición**.

La Sociedad Humana

Falsas Expectativas sobre la IA Gracias al “Efecto Hollywood”



The Terminator (Wikipedia).

- Sí, la IA está cambiando las cosas.
- Pero hasta el momento sólo hay **IAs angostas**: se necesitan a las personas antes, durante y después de la introducción de estas tecnologías.
- La tecnología actual está **muy lejos de ser capaz de duplicar la enorme capacidad humana**.
- Queda por verse **hasta dónde y cuán rápido** va a ser el cambio.

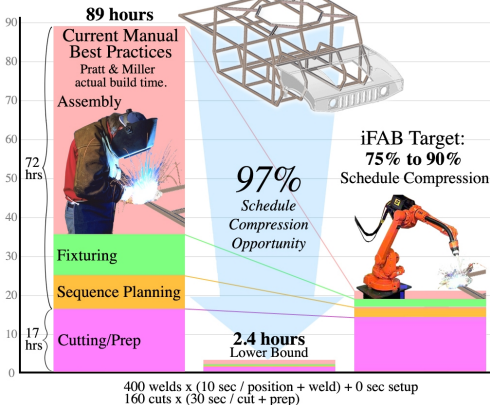
La Sociedad Humana

Predicciones sobre la Automatización del Trabajo

When	Where	Jobs Destroyed	Jobs Created	Predictor
2016	worldwide		900,000 to 1,500,000	Metra Martech
2018	US jobs	13,852,530*	3,078,340*	Forrester
2020	worldwide		1,000,000-2,000,000	Metra Martech
2020	worldwide	1,800,000	2,300,000	Gartner
2020	sampling of 15 countries	7,100,000	2,000,000	WEF
2021	worldwide		1,900,000-3,500,000	IFR
2021	US jobs	9,108,900*		Forrester
2022	worldwide	1,000,000,000		Thomas Frey
2025	US jobs	24,186,240*	13,604,760*	Forrester
2025	US jobs	3,400,000		ScienceAlert
2027	US jobs	24,700,000	14,900,000	Forrester
2030	worldwide	2,000,000,000		Thomas Frey
2030	worldwide	400,000,000-800,000,000	555,000,000-890,000,000	McKinsey
2030	US jobs	58,164,320*		PWC
2035	US jobs	80,000,000		Bank of England
2035	UK jobs	15,000,000		Bank of England
No Date	US jobs	13,594,320*		OECD
No Date	UK jobs	13,700,000		IPPR

Every study we could find on what automation will do to jobs, in one chart, Erin Winick, Technology Review, 25 de enero de 2018.

iFab Opportunity



Instruction Generation

Robotic Assembly

MACHINE INSTRUCTION



Augmented Fixturing

HUMAN INSTRUCTION



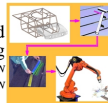
Virtual Layout

HUMAN INSTRUCTION



Automated Process Planning

MACHINE & HUMAN INSTRUCTION



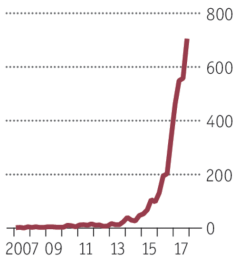
Robots and People Can Work Faster Together, David Bourne, Director del Rapid Manufacturing Lab, Robotics Institute, Carnegie Mellon University, 25 de julio de 2013. Leer *Collaborative Intelligence: Humans and AI Are Joining Forces*, Harvard Business Review, Julio-Agosto, 2018.

Adaptaciones de las Empresas

Reacciones

Machine earning

Mentions of AI and machine learning on earnings calls of public companies



Source: Bloomberg

Economist.com

Non-tech businesses are beginning to use artificial intelligence at scale, The Economist, 31 de marzo de 2018.

- ▶ Mercados horizontales compuestos por consumidores generales están dominados por **Alibaba, Amazon, Apple, Baidu, Facebook, Google, IBM, Microsoft y Tencent.**
- ▶ El resto son **silos aislados.**
- ▶ **Quién tiene la gente y los datos domina.**
- ▶ Los algoritmos se están transformando en un **commodity.**

Adaptaciones de las Empresas

IA en la Internet



Adaptaciones de las Empresas

IA en la Empresa



AI Lanscape 2018, Topbots, septiembre de 2018.



Adaptaciones de las Empresas

IA en el Mundo Físico

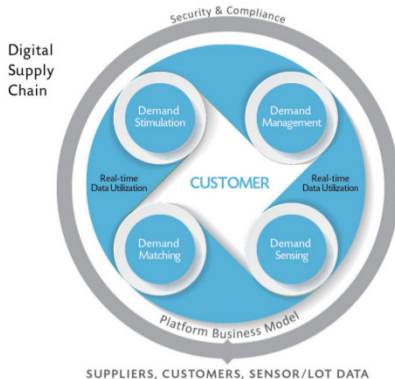


YuMi de ABB.

- ▶ **El último paso.**

Adaptaciones de las Empresas

Ciclo de Percepción y Acción



DRIVING DEMAND IN THE DIGITAL SUPPLY CHAIN: Algorithms and the Untapped Power of Applying Real-Time Big Data and AI/ML, DSCI Institute, The Center for Global Enterprise, enero de 2018.

Adaptaciones de las Empresas

Data Driven

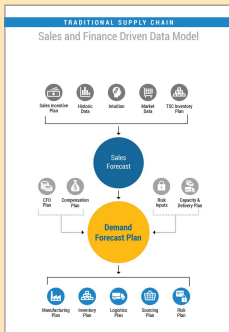


EXHIBIT 1

Driving Demand in the Digital Supply Chain

38

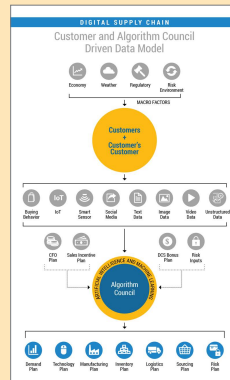


EXHIBIT 2

Driving Demand in the Digital Supply Chain

39

DRIVING DEMAND IN THE DIGITAL SUPPLY CHAIN: Algorithms and the Untapped Power of Applying Real-Time Big Data and AI/ML, DSCI Institute, The Center for Global Enterprise, enero de 2018.

- ▶ Las empresas en el mundo están **acelerando** la introducción de la inteligencia artificial.
- ▶ Es necesario educar a las personas, **partiendo por los líderes**.
- ▶ Fuerte cambio de cultura interna que **rompe los silos** típicos.
- ▶ Es necesaria la ingeniería **antes, durante y después**.

Conclusión

¿Preguntas?