**START PAGE**

MARIE SKLODOWSKA-CURIE ACTIONS

**Individual Fellowships (IF)**
**Call: H2020-MSCA-IF-2015**

PART B

"OSEGA"

**This proposal is to be evaluated as:**

**[Standard EF]**

# Contents

## 0  List of Participants

| Participants | Legal Entity Short Name | Academic | Non-academic | Country | Dept. / Division / Laboratory | Supervisor | Role of Partner Organisation |
|---|---|---|---|---|---|---|---|
| Beneficiary | | | | | | | |
| Lancaster University | ULANC | X | | United Kingdom | Department of Mathematics and Statistics | Prof David Leslie | |

# 1 Excellence

## 1.1 Quality, innovative aspects and credibility of the research

A critical concern in the modern world is security. Effectively protecting ports, airports, trains and other transportation systems from malicious attacks, combating the trafficking of drugs, firearms and even people, and securing proprietary and sensitive information over the ever-growing cyber-networks, comprise some of the principal axes of this critical task. The main challenge in all of these problems is that maximum security must be obtained with a limited number of available resources. For instance, the total number of security agents is typically less than the number of targets that need to be protected. This calls for the design of appropriate resource allocation techniques which optimise security under the constrained resources available, in the presence of uncertainty about the adversaries' interests.

Security resource allocation and scheduling problems comprise one of the many application areas that have recently been shown to greatly benefit from game-theoretic approaches. Indeed, as a solid mathematical framework to model strategic decision making, game theory has proved useful in many real-world applications from economics and political science to logic, computer science and psychology. In this paradigm, the problem is cast as a "game" and the objective is to find a solution whereby each "player" makes choices to maximise her own *utilities*, which may often be in conflict with those of her opponent. A "security game" corresponds to a competition between a defender and an attacker. To solve a security game, all possible actions (attacks and defences) of the two players are enumerated, and for each player an outcome (value) is assigned, which depends on the pair of actions taken by both players. In cases where these outcomes are known, game-theoretic approaches have provided impressive results. Since 2007, the so-called ARMOR software[1] has been used at Los Angeles International Airport to effectively determine checkpoints on roadways leading to the airport, and to determine canine patrol routes within terminals. Similar deployments have been made by the US Federal Air Marshals,[2] the US coast guard,[3] and to design the Los Angeles Metro system's fare inspection strategy.[4]

A severe limitation of these models is that they assume the utility functions to be known, whereas they must actually be estimated by experts or obtained from historical data. As a result, potentially high estimation errors or a lack of historical data (which is inevitable in a quickly-evolving security scenario) may render the security game solver useless. Therefore it is of importance to use methods that can quickly collect the most relevant data in order to estimate the parameters of the game and quickly reach a satisfactory operational performance.

This project will therefore deploy approaches of *statistics* and *machine learning* within the framework of security games which are played repeatedly between a defender and attackers. Repeated security games allow for the continuous collection of data, which can in turn be used to estimate the parameters of the game and influence future behaviour. Our key objective is to design efficient and theoretically sound, data-driven methods that can actively interact with the environment to *learn* and *act* in security games of realistic scales, which must therefore be *practical*, *scalable* and *robust*.

**State-of-the-art**

**Security meets Game Theory**  From a game-theoretic perspective, a security problem is viewed as a two-player game that captures the interaction between a defender (e.g., border patrols, metro inspectors, network administrators) and attackers (e.g., terrorists/smugglers, illegal metro users, malicious cyber attackers). The action of the defender (attacker) is defined as selecting a subset of targets to protect (attack). For each defender/attacker action pair, *utilities* are defined as the players' gain or loss, and the players'

---

[1]James Pita et al. "Deployed ARMOR protection: the application of a game theoretic model for security at the Los Angeles International Airport". In: *Proceedings of the 7th international joint conference on Autonomous agents and multiagent systems: industrial track*. International Foundation for Autonomous Agents and Multiagent Systems. 2008, pp. 125–132.

[2]Jason Tsai et al. "IRIS-a tool for strategic security allocation in transportation networks". In: (2009).

[3]Eric Shieh et al. "Protect: A deployed game theoretic system to protect the ports of the united states". In: *Proceedings of the 11th International Conference on Autonomous Agents and Multiagent Systems-Volume 1*. International Foundation for Autonomous Agents and Multiagent Systems. 2012, pp. 13–20.

[4]Zhengyu Yin et al. "TRUSTS: Scheduling randomized patrols for fare inspection in transit systems using game theory". In: *AI Magazine* 33.4 (2012), p. 59.

objectives are to maximise their corresponding pay-offs. From the defender's perspective, this corresponds to efficiently allocating a limited number of resources to secure some predefined targets. The expected utilities for both players can be stored in two matrices, $\boldsymbol{A}$ for the defender, $\boldsymbol{B}$ for the attacker. Entries $\boldsymbol{A}_{i,j}$ and $\boldsymbol{B}_{i,j}$ are the expected utilities for the defender and attacker, respectively, when the defender plays strategy $i$ and the attacker responds with strategy $j$. Solutions to such games rely on randomised (mixed) strategies, making each player's behaviour unpredictable to the other. If a fully competitive setting, in which the attacker's gain is the defender's loss (i.e. the so-called zero-sum games) a fully robust strategy can be calculated for the defender, in that it provides guaranteed performance against *any* possible attacker, even if the defender's strategy is completely revealed to the attacker. A generalisation to more general (non zero-sum) games, which is particularly relevant to security games, is the Stackelberg equilibrium[5] in which the defender's mixed strategy is first publicised and the attacker plays a best response to this mixed strategy. It is this Stackelberg security game framework that will be considered in the proposed project.

**Uncertainty in Security Games** Most standard game-theoretical analyses assume that the payoff matrices are known in advance. However uncertainty is endemic in most real-world applications. For instance, in the context of security games, the random selection of passengers for security checks at an airport is a source of uncertainty in this game, where the outcome is random and the probability of successful security enforcement is unknown. As also confirmed by several empirical studies in fraud and cybercrime detection, this phenomenon can significantly decrease the defender's performance.[6] Extensive studies have been dedicated to the design of security games that are robust with respect to uncertainty about the environment.[7] However, an important observation is that much more can be done in the case of *repeated* security games. Indeed, this repetition allows the defender to further reduce her uncertainty about the model and intelligently *learn* how to improve her performance over time. Specifically, in this case, a security game solver can autonomously take intelligent decisions at repeated instances of the game, by carefully collecting, extracting and acting upon information from historical data. As discussed further in below, this is precisely the area where mathematical machine learning has the strongest results.

**Bandit problems** Mathematical machine learning is a modern amalgamation of statistics (to deal with uncertain data) and optimisation (to efficiently select appropriate actions). One of the fundamental problems in machine learning, relevant to our research objectives in this proposal, is the *multi-armed bandit problem*. It corresponds to a scenario where the learner is required to actively collect data from an environment in order to solve a given task. Solutions built for this problem have found many practical applications from adaptive routing in a network to medical trials of new medicines.[8] A multi-armed bandit is precisely a game as described above, but where the attacker has only one action; the bandit problem results when the game is repeated, and on iteration $t$ the learner selects action $i(t)$ and received a reward $l_t$ which is a random variable with expectation equal to $\boldsymbol{A}_{i_t,1}$. An important constraint in this setting is that the player is not allowed to observe the hypothetical reward that would have been collected had another arm been selected instead.

In the bandit problem two related tasks are commonly considered. One is the *online decision-making* task, in which the reward for each decision must be taken into account; this task is relevant to the immediate

[5]Dmytro Korzhyk et al. "Stackelberg vs. Nash in Security Games: An Extended Investigation of Interchangeability, Equivalence, and Uniqueness." In: *J. Artif. Intell. Res.(JAIR)* 41 (2011), pp. 297–327.

[6]Jennifer S Granick. "Faking It: Calculating Loss in Computer Crime Sentencing". In: *ISJLP* 2 (2005), p. 207; Peter Swire. "No cop on the beat: Underenforcement in e-commerce and cybercrime". In: *J. on Telecomm. & High Tech. L.* 7 (2009), p. 107.

[7]Michele Aghassi and Dimitris Bertsimas. "Robust game theory". In: *Mathematical Programming* 107.1-2 (2006), pp. 231–273; Thanh H Nguyen et al. "Regret-based optimization and preference elicitation for stackelberg security games with uncertainty". In: *Proceedings of the 28th AAAI Conference on Artificial Intelligence (AAAI)*. 2014, pp. 756–762; Christopher Kiekintveld, Towhidul Islam, and Vladik Kreinovich. "Security Games with Interval Uncertainty". In: *Proceedings of the 2013 International Conference on Autonomous Agents and Multi-agent Systems*. AAMAS '13. St. Paul, MN, USA: International Foundation for Autonomous Agents and Multiagent Systems, 2013, pp. 231–238. ISBN: 978-1-4503-1993-5.

[8]Sébastien Bubeck and Nicolo Cesa-Bianchi. "Regret analysis of stochastic and nonstochastic multi-armed bandit problems". In: *arXiv preprint arXiv:1204.5721* (2012).

deployment of the system to actually make decisions while it learns. The standard performance metric here is the regret, defined to be the difference between the total reward that could have been achieved if full information were available in advance, and the actual reward that was achieved. The other task is that of *pure exploration*, in which a learning phase is permitted during which received rewards do not matter; this corresponds to allowing a training phase for the system prior to deployment. The performance metric of this problem is the quality of the arm selected immediately after the training phase. Application of the pure-exploration problem to parallel action selection in robotic planning has been extensively studied by Dr Gabillon.[9]

When we consider more general security games, in which the attacker has more than one available action, we can instead think of a non-strategic attacker who uses a fixed distribution over actions through time. This also corresponds to bandit problem for the attacker. Bandit strategies are therefore important for active learning in security games. However it is as yet an important open question how to take into account strategic or adaptive behaviour of the attacker. It is this aspect of learning in security games which will be developed in the proposed project.

**Existing results**  Recent advances have been made in this direction, which mostly focus on the case where the attacker's preferences are not fully known and must be learned through repeated plays of the game. One approach is to analyse the number of required queries to learn the optimal defender's strategy.[10] Alternatively a Bayesian approach can be taken, with techniques based on Partially Observable Markov Decision Processes (POMDPs) used to update a posterior over the adversary's preferences.[11] Recently, a more relevant analysis has been given[12] for the case of multiple attackers, where at each round of the game, a single attacker is chosen adversarially from a fixed, finite, set of known attackers. This corresponds to a case where the utility matrix $\boldsymbol{B}$ is chosen adversarially from a set of $k$ known matrices, and shows strong connections with adversarial bandit theory. However all of these security games results rely restrictive assumptions about prior knowledge and observability, and on the number of available actions to each player being reasonably small. In this proposal we consider combinatorial actions (selecting $k$ resources to protect out of $n$ that may be attacked) and so the total number of actions available is enormous.

**Vision**  The purpose of this project is to create *practical*, *scalable* and *robust* methods for security games. First, we target *practicality* in the sense that our methods will be autonomous in handling uncertainty in the model and will actively reduce this uncertainty by interacting with the environment in which the game takes place during repeated plays of the game. We will do so by combining existing security game research with the extremely active field of multi-armed bandit research. Second, we target *scalability* so that the methods will apply with an extremely large number of possible actions. This can be achieved by making simplifying structure assumptions such as combinatorial structure, in which an action consists of selecting $k$ objects from $n$. Finally *robustness* is a key issue in security games in several senses. Methods should not break down in the face of adversarial play by the attacker both during learning (so that the defender receives the least helpful information from which to learn). Furthermore methods should not perform well only in average — worst case performance is extremely important. We will therefore devise a theoretically

---

[9]V. Gabillon et al. "Multi-Bandit Best Arm Identification". In: *Proceedings of the Advances in Neural Information Processing Systems 25*. 2011, pp. 2222–2230.

[10]Avrim Blum, Nika Haghtalab, and Ariel D Procaccia. "Learning optimal commitment to overcome insecurity". In: *Advances in Neural Information Processing Systems*. 2014, pp. 1826–1834; Joshua Letchford, Vincent Conitzer, and Kamesh Munagala. "Learning and approximating the optimal strategy to commit to". In: *Algorithmic Game Theory*. Springer, 2009, pp. 250–262.

[11]Janusz Marecki, Gerry Tesauro, and Richard Segal. "Playing Repeated Stackelberg Games with Unknown Opponents". In: *Proceedings of the 11th International Conference on Autonomous Agents and Multiagent Systems - Volume 2*. AAMAS '12. Valencia, Spain: International Foundation for Autonomous Agents and Multiagent Systems, 2012, pp. 821–828. ISBN: 0-9817381-2-5, 978-0-9817381-2-3; Yundi Qian et al. "Online planning for optimal protector strategies in resource conservation games". In: *Proceedings of the 2014 international conference on Autonomous agents and multi-agent systems*. International Foundation for Autonomous Agents and Multiagent Systems. 2014, pp. 733–740.

[12]Maria-Florina Balcan et al. "Commitment without regrets: Online learning in Stackelberg security games". In: (2015).

sound approach in which efficient algorithms are developed and finite time performance guarantees are provided.

**Objective 1: Scalability in Pure exploration Bandits via Submodularity.** When considering the combinatorial nature of many security games, it is necessary to use this structure intelligently. In particular, naively enumerating all possible actions as in the standard formulation of security games makes the computations rapidly intractable. We will address this problem by first studying combinatorial bandit techniques, which will subsequently be extended to use in security games. Combinatorial bandits form a central part of Dr Gabillon's area of expertise. A key observation is that in many cases the players' performance utility function is submodular. One example is in the context of maximal coverage problem for sensor (or checkpoint) placement.[13] This submodularity property can in turn be used to provide tractable and almost optimal algorithms in the online decision-making settings of bandits.[14] However pure exploration has not yet been solved in this area; Objective 1 is therefore to complete the analysis of combinatorial bandits by addressing the pure exploration problem under submodularity assumptions.

**Objective 2: Pure Exploration in Security Games.** As discussed previously, a major barrier to applying security games to real-world scenarios is that the players' utility matrices $A$ and $B$ are unknown and must be learned. In the first of two objectives on standard security games we consider the case where the defender is in fact able to safely examine her defensive strategies before applying them online, corresponding to the pure exploration bandit framework. A familiar real-world example is the security system at an airport, which can be tested many times before being deployed as the principal defence scheme. In this formulation, we assume both utility matrices $A$ and $B$ unknown. During the test phase, the defender is able to probe an entry of its utility matrix at every mock repetition of the game. The value obtained is *a noisy version* of the true entry $A_{i,j}$.

We will design a strategy for the defender to either minimise the number of tests needed to identify an excellent strategy with a given level of confidence[15] or to maximise her probability of identifying the best strategy given a fixed number of tests.[16] In order to extend these classical results to the setting proposed above, we will first carefully characterise the data-dependent hardness of the problem, extending recent relevant results for combinatorial bandits[17] and similar results currently in preparation by Dr Gabillon. Of course, in games these complexity results are more challenging than in the bandit problem since the complexities depend also on the actions available to the attacker. Therefore we will study this problem by gradually increasing its difficulty with different partial feedback structures. First we note that a recent work[18] on query complexity, corresponds to the simpler deterministic version of this problem where it is assumed that the probing outcome corresponds to the *true value* of $A_{i,j}$. Therefore our first approach will be to combine ideas from pure exploration and the deterministic query complexity setting in a context where the defender can individually sample from any entry of the matrix. A second more challenging setting will be to consider the more adversarial learning problem where the defender chooses a strategy and only observes the value of the game corresponding to an action selected by the attacker; the attacker might be either oblivious to the defender, or playing a Stackelberg best-response to the defender strategy.

**Objective 3: Regret Analysis of Repeated Security Games.** In some applications a newly

---

[13] Andreas Krause, Alex Roper, and Daniel Golovin. "Randomized sensing in adversarial environments". In: *IJCAI Proceedings-International Joint Conference on Artificial Intelligence*. Vol. 22. 3. 2011, p. 2133.

[14] Victor Gabillon et al. "Adaptive submodular maximization in bandit setting". In: *Advances in Neural Information Processing Systems*. 2013, pp. 2697–2705.

[15] O. Maron and A. Moore. "Hoeffding races: Accelerating model selection search for classification and function approximation". In: *Proceedings of the Advances in Neural Information Processing Systems 7*. 1993; E. Even-Dar, S. Mannor, and Y. Mansour. "Action Elimination and Stopping Conditions for the Multi-Armed Bandit and Reinforcement Learning Problems". In: *Journal of Machine Learning Research* 7 (2006), pp. 1079–1105.

[16] S. Bubeck, R. Munos, and G. Stoltz. "Pure Exploration in Multi-Armed Bandit Problems". In: *Proceedings of the Twentieth International Conference on Algorithmic Learning Theory*. 2009, pp. 23–37; J.-Y. Audibert, S. Bubeck, and R. Munos. "Best Arm Identification in Multi-Armed Bandits". In: *Proceedings of the Twenty-Third Conference on Learning Theory*. 2010, pp. 41–53.

[17] Shouyuan Chen et al. "Combinatorial pure exploration of multi-armed bandits". In: *Advances in Neural Information Processing Systems*. 2014, pp. 379–387.

[18] Paul W Goldberg and Stefano Turchetta. "Query Complexity of Approximate Equilibria in Anonymous Games". In: *arXiv preprint arXiv:1412.6455* (2014).

created security system is not provided with any historical data and cannot be tested before being used in production. Here, the learning of the utilities must be performed online while actually playing the security game. In this online decision-making context it is of high importance for the agent to learn the utilities as fast as possible. This means that only providing an analysis to demonstrate asymptotic convergence[19] is inadequate. To address repetitive learning in security games we will here assume that the utility matrix $A$ is unknown to the defender and that, at each repetition of the game the attacker will best respond to the defender's current strategy (if $\pi$ is a mixed strategy of the defender, then denote $b(\pi)$ the best response of the attacker). Therefore we are considering a setting that is related to the analysis of Stackelberg equilibrium. A natural quantity of interest is the *cumulative regret*, defined as

$$R(n) = n \max_{\pi} [\pi A b(\pi)] - \sum_{t=1}^{n} \pi_t A b(\pi_t).$$

This compares the actual reward received (assuming the attacker always performs as well as they can) to the reward achieved at Stackelberg equilibrium (when the defender chooses the best possible mixed strategy under the knowledge that the attacker will best respond to it). The objective is for the defender to build a series of defence strategies $\pi_t$ for $t = 1, \ldots, n$ to minimize the expected cumulative regret $R(n)$.

This game-theoretic scenario, is actually strongly related to the bandit problem, in that the best-responding assumption on the attacker leads to a situation where the reward to the defender depends only on the selected (mixed) strategy. A solution can thus be obtained by considering it as a bandit problem with continuous action space consisting of the set of all probability distributions on the original discrete action space. However a more efficient solution is likely to be obtained by explicitly considering the game-theoretical nature. In particular, most current approaches for online decision-making in bandits, such as upper confidence bound methods,[20] implement a strategy that is optimistic in face of uncertainty, playing any strategy that could be the best given the level of uncertainty. It will be extremely interesting to discover if this optimism principle still holds in an adversarial game, or whether a more cautious approach is needed. Approaches based on Thompson sampling[21] will be also considered as they have proved very efficient in practice and correspond to an area of expertise of Prof. Leslie.

Finally an interesting additional requirement is to learn security strategies that are not only of good quality in average but also whose performance is not subject to large variance when used on a daily basis. This *risk-averse* requirement has been well-studied in the statistical community, and has recently been considered in a multi-armed bandit framework.[22] An implementation in the security games specific context is a very natural extension to the main body of work in this objective.

**Objective 4: Learning in combinatorial games.** Objective 4 will be devoted to solving security games with more complex action structures. Real-world security problems often involve large, complex networks. This includes, for instance, complex routes, or computer/communication networks. The size of the action spaces for both attacker and defender is often combinatorially large. Standard results with convergence times that increase with the size of the action spaces become extremely weak in such settings. We will therefore take advantage of the inherent combinatorial structure of the problem to create efficient and computationally tractable algorithms in these large games. In particular, we will develop the approaches of both Balcan[23] in simple Stackelberg games and in adversarial combinatorial bandits.[24] The issue of

---

[19]David S. Leslie and E.J. Collins. "Generalised weakened fictitious play". In: *Games and Economic Behavior* 56.2 (2006), pp. 285–298. ISSN: 08998256. DOI: 10.1016/j.geb.2005.08.005. URL: http://linkinghub.elsevier.com/retrieve/pii/S089982560500103X; Archie C. Chapman et al. "Convergent Learning Algorithms for Unknown Reward Games". en. In: *SIAM Journal on Control and Optimization* 51.4 (Jan. 2013), pp. 3154–3180. ISSN: 0363-0129. DOI: 10.1137/120893501. URL: http://epubs.siam.org/doi/abs/10.1137/120893501.

[20]P. Auer, N. Cesa-Bianchi, and P. Fischer. "Finite-Time Analysis of the Multi-Armed Bandit Problem". In: *Machine Learning* 47 (2002), pp. 235–256.

[21]Daniel Russo and Benjamin Van Roy. "An information-theoretic analysis of thompson sampling". In: *arXiv preprint arXiv:1403.5341* (2014).

[22]Amir Sani, Alessandro Lazaric, and Rémi Munos. "Risk-Aversion in Multi-armed Bandits". In: *Advances in Neural Information Processing Systems 25*. Ed. by F. Pereira et al. Curran Associates, Inc., 2012, pp. 3275–3283.

[23]Balcan et al., "Commitment without regrets: Online learning in Stackelberg security games".

[24]Nicolo Cesa-Bianchi and Gábor Lugosi. "Combinatorial bandits". In: *Journal of Computer and System Sciences* 78.5 (2012),

scalability will be addressed in light of the results found in Objective 1.

**Objective 5: Repeated Network-Security Games.** As a more concrete application of Objective 4, this objective will focus on the particular combinatorial structure that is a graph as this stucture is present in numerous real-word applications. In light of the ever-growing, modern, social and communication networks, a canonical example is that of smuggler arrest in a network.[25] This has received significant attention in the community, especially in response to the Mumbai attacks of 2008, after which Mumbai Police started to schedule a limited number of inspection checkpoints on the road throughout the city. This problem has not yet been studied in its repeated form, where a pursuit-evasion game is played multiple times against a population of smugglers. Therefore, the currently deployed strategy of the defender is *not adaptive* to observations collected about the attackers' historical strategy and is therefore sub-optimal. We will therefore develop adaptive strategies, using the approaches developed in Objectives 1–4. However, the graph structure provides additional constraints on the action spaces, and provides additional information, when compared with the general combinatorial problem. In particular, the set of actions available to the attacker is restricted to a set of paths through the network, and the the graph structure provides strong information about sensible choices of checkpoints (for example, aligning them all along one route through the network is a particularly unfortunate choice, but is not ruled out by a generic combinatorial structure). Furthermore, absence or presence of a smuggler on one link of the graph will likely provide information about which other other links were utilised on that iteration. Therefore the objective is to design specific algorithms in in situations where it is possible to take advantage of specific graphical structure of the problem. We will start by defining a notion that captures the hardness of the task depending on characteristics of the graph that we would discover. Note that, although the goal is to generate algorithms for security games on graphs, the results to be obtained will be expected to lay grounds for research in a more general setting of active learning with graph structure. Dr Gabillon has held initial discussions on this topic with Dr Michal Valko, a world-famous expert in active learning on graphs and part of INRIA Lille in France. This project will allow the formation of a productive and lasting collaboration with Dr Valko, which will be greatly beneficial not only in achieving the this objective, but also to strengthen international links between Lancaster University in the UK and INRIA in France.

## 1.2 Clarity and quality of transfer of knowledge/training for the development of the researcher in light of the research objectives

The overall training objective is to significantly develop Dr Gabillon's scientific, organisational, communication and technology transfer skills. This will enable him to continue building his portfolio of outstanding research to attain a position of independence and gain recognition in the international research community.

The proposed project is primarily a research project, and the main training objectives are to enhance the Dr Gabillon's scientific skills. Dr Gabillon is already an expert in the modern theory of bandits, including best arm identification, and reinforcement learning. Therefore this project's main training objective will be to develop his skills and knowledge in statistical learning methods and game theory. Prof. Leslie is an expert in both areas, and will of course assist the development of Dr Gabillon. Further expertise in Lancaster from whom Dr Gabillon will learn includes the Statistical Learning group, in which he will be based, and the broader Statistics Research Group.

Combining the Department of Mathematics and Statistics with the Operations Research group within the Management School, Lancaster is the leading UK institution in bandit theory, with expertise in index policies (Glazebrook, Kirkbride, Jacko), Thompson sampling and contextual bandits (Grunewalder, Leslie) and application in medical trials (Vilar). Dr Gabillon will have ample opportunity to further develop his expertise in this area, and indeed brings expertise from a complementary aspect of online learning and decision-making in the design and analysis of algorithmic approaches to learning, especially with combinatorial bandit problems. Dr Gabillon's expertise in best-arm identification will be of great interest

---

pp. 1404–1422.

[25]Manish Jain et al. "A double oracle algorithm for zero-sum security games on graphs". In: *The 10th International Conference on Autonomous Agents and Multiagent Systems-Volume 1*. International Foundation for Autonomous Agents and Multiagent Systems. 2011, pp. 327–334.

to the Medical and Pharmaceutical Statistics research group. He will present his research in this area to the research group and discuss possible applications in clinical trial design. Furthermore his expertise in combinatorial bandits complements current industrially-funded research of the Prof. Leslie.

In addition to his research skills, Dr Gabillon will learn from the host's world-leading expertise in developing industrially-inspired statistics. Statistical researchers in Lancaster have constant exposure to external companies, through the STOR-i Centre for Doctoral Training, and the Data Science Institute. While embedded in this culture, Dr Gabillon will be given the opportunity to:

1. Gain further experience of developing industry/academic partnerships by working with Profs. Leslie and Eckley and other staff in STOR-i and the Data Science Institute in technology transfer activities.

2. Develop public communication skills by presenting research results to varied audiences.

3. Participate in the organisation of workshops in Lancaster and at the Royal Statistical Society.

4. Receive training on applying for funding by co-authoring proposals for UK and EU funding agencies.

5. Attend staff training workshops designed specifically for early-career researchers, and specifically the Research Development Programme, a structured development route for researchers, designed to promote impactful research and to support development beyond a disciplinary area.

6. Participate in teaching and research supervision (undergraduate and graduate). This will not be obligatory, but Dr Gabillon will have the opportunity to benefit from peer observation and mentoring.

The UK Concordat to Support the Career Development of Researchers is an agreement between funders and employers of research staff to improve the employment and support for researchers and research careers in UK higher education. Lancaster University is fully committed to the Concordat to Support the Career Development of Researchers and has put in place an Action Plan to support the full implementation of the Concordat at Lancaster. Furthermore, throughout the fellowship, Dr Gabillon will adhere to the European Charter for Researchers, and the training objectives will be managed through a Personal Career Development Plan that Prof. Leslie and Dr Gabillon will write together. This plan will be revised regularly throughout the fellowship to ensure that all objectives are met. In addition, Dr Gabillon will have regular meetings with Prof. Leslie to discuss his research and to receive advice.

## 1.3 Quality of the supervision and the hosting arrangements
**Qualifications and experience of the supervisor(s)**

Prof. Leslie leads the Statistical Learning research group in the Department of Mathematics and Statistics, Lancaster University, and is Theme Lead for Foundations in Lancaster University's new Data Science Institute. He is a world-leading researcher in statistical learning, Bayesian inference, decision-making and game theory, with 19 refereed articles in top journals of several different research fields, and collaborators from France, Singapore, USA and Australia. His research on contextual bandit algorithms[26] is used by many of the world's largest companies to balance exploration and exploitation in real-time website optimisation. He is expert in the mathematics of learning in games,[27] stochastic approximation,[28] and the

---

[26]BC May et al. "Optimistic Bayesian sampling in contextual-bandit problems". In: *The Journal of Machine Learning Research* 13 (2012), pp. 2069–2106. URL: http://dl.acm.org/citation.cfm?id=2343711.

[27]David S Leslie and E J Collins. "Convergent Multiple-timescales Reinforcement Learning Algorithms in Normal Form Games". In: *Annals of Applied Probability* 13.4 (2003), pp. 1231–1251; David S Leslie and E J Collins. "Individual Q-learning in normal form games". In: *SIAM Journal on Control and Optimization* 44 (2005), pp. 495–514; Leslie and Collins, "Generalised weakened fictitious play"; Chapman et al., "Convergent Learning Algorithms for Unknown Reward Games"; Steven Perkins and David S. Leslie. "Stochastic Fictitious Play with Continuous Action Sets". In: *Journal of Economic Theory* 152 (2014), pp. 179–213. DOI: 10.1016/j.jet.2014.04.008. URL: http://www.sciencedirect.com/science/journal/00220531.

[28]Leslie and Collins, "Convergent Multiple-timescales Reinforcement Learning Algorithms in Normal Form Games"; Steven Perkins and David Leslie. "Asynchronous stochastic approximation with differential inclusions". en. In: *Stochastic Systems* 2 (2012), pp. 409–446. DOI: 10.1214/11-SSY056.. URL: http://www.i-journals.org/ssy/viewarticle.php?id=56; Perkins and Leslie, "Stochastic Fictitious Play with Continuous Action Sets".

mathematics of statistically-inspired reinforcement learning.[29] Prof. Leslie is the holder of a Google Faculty Award which funds a student to investigate multiple-action selection in bandits. Prior to his relocation to Lancaster, he was a senior lecturer in the statistics group of the School of Mathematics, University of Bristol. He continues to be co-director of the £1.5m EPSRC-funded cross-disciplinary decision-making research group at the University of Bristol, and was on the management team of the £5.5m ALADDIN project, a large strategic partnership between BAE Systems and EPSRC, involving researchers from Imperial College, Southampton, Oxford, Bristol and BAE Systems.

Prof. Leslie's mentoring approach is one of 'guided freedom' in which the mentee takes responsibility for their own research, while regular discussions ensure that dead ends are avoided and promising openings are exploited. In the 10 years since taking up a Faculty position, he has supervised 17 PhD students (5 now in Academic positions), 2 post-doctoral fellows, numerous MSc and undergraduate dissertations, and an undergraduate secondment from ENS Lyon.

**Hosting arrangements**

The Researcher will be embedded within the statistical learning group which is lead by Prof. Leslie. This is a team of 5 academic staff and around 5 PhD students within the Department of Mathematics and Statistics. Dr Gabillon will participate in weekly group meetings and benefit from advice from the senior scientists in the group on research direction and management, personal development, workshop organisation, teaching, and other aspects of academic life. The group also has extremely strong links with both the Data Science Institute (www.lancaster.ac.uk/dsi/) and the STOR-i Centre for Doctoral Training (www.stor-i.lancs.ac.uk/); each provides a weekly seminar series. These exciting initiatives will provide multiple further opportunities to develop informal mentoring relationships in addition to the formal process which takes place for all staff at Lancaster University. To ensure integration within these networks Dr Gabillon will be introduced to the groupings of researchers, invited to deliver a seminar on his research, and will participate in away days in which strong relationships are developed.

## 1.4 Capacity of the researcher to reach and re-enforce a position of professional maturity in research

Professional maturity in academia would be ideally reached by leading a dynamic research group actively working on fundamental problems at the interface of game theory and online learning, with strong impact in real-world applications. Dr Gabillon has shown an extremely high potential to achieve this goal. As evident from his strong publication record, he has broad expertise in the domain, always giving equal importance to theory and applications. Dr Gabillon has also demonstrated strong ability to acquire new knowledge and become highly productive in a short period of time. Indeed, a significant result of his PhD thesis is in bringing classical reinforcement learning algorithms closer to daily life. During the course of his PhD, through a 6-months internship at a major US R&D lab (Technicolor Research Laboratory, Palo Alto), he had the opportunity to collaborate with a new team of R&D researchers. He quickly became productive and his efforts in this short period of time have resulted in the publication of two peer-reviewed papers at prestigious international conferences in machine learning. Through this experience he has also obtained valuable knowledge about industrial research and its interaction with academia. This has given him the ability to better understand the research pathways to produce high-impact results and establish significant collaborations with industry.

At the start of the fellowship, Dr Gabillon will be closely mentored by Prof. Leslie at Lancaster University. He will also have access to the university's research resources, and will be able to further develop his research and supervision skills, which will greatly contribute to achieving professional maturity. At Lancaster University, Dr Gabillon will also have the unique opportunity to establish inter-disciplinary collaborations through the recently established STOR-i program, a quality research training interface between statistics and industry.

---

[29]Leslie and Collins, "Individual Q-learning in normal form games"; Tobias Larsen et al. "Posterior Weighted Reinforcement Learning with State Uncertainty". In: *Neural Computation* 22 (2010), pp. 1149–1179.

## 2  Impact

### 2.1  Enhancing research- and innovation-related skills and working conditions to realise the potential of individuals and to provide new career perspectives

Dr Gabillon is already a leading researcher in the mathematics of bandit algorithms and reinforcement learning. This fellowship provides a training opportunity in two key additional research competences. Firstly, he will develop an in depth knowledge of cutting edge statistical theory, and bring that to bear within bandit algorithms. Training will be received from leading scientists in statistics and operations research at Lancaster University, and the many international visiting researchers who visit the department. Secondly, Prof. Leslie is a leading expert on learning in games, as well as bandit algorithms, and will mentor Dr Gabillon to bring ideas from bandits into the game theoretical scenarios of this research proposal. This significant broadening of the researcher's skill set will give him an extremely solid foundation on which to build a future research career.

In addition to pure research opportunities, Dr Gabillon will work within Lancaster University's extremely effective framework for industrial collaboration. He will develop skills in how to manage the industry/academia relationship to ensure mutually beneficial outcomes. This relationship-management will be a key skill for academics in the future; Lancaster University, and particularly the Department of Mathematics and Statistics, is currently a world-leading institution in developing such relationships. Dr Gabillon will both be introduced to prospective industrial partners, and receive mentoring as he develops his own relationships.

### 2.2  Effectiveness of the proposed measures for communication and results dissemination

With the launch of the Data Science Institute, Lancaster University will be inaugurating a "Data Science Network", in conjunction with Lancaster University's Knowledge Business Centre, an innovation hub providing a gateway for business/academic interaction which allows the transfer of expertise between Lancaster's academics, regional businesses and community partnerships through training and technology transfer activities. This network will bring together academic data scientists with local companies in regular show and tell sessions. Dr Gabillon will be a regular participant at these events, enabling bi-directional communication of opportunities and requirements, and the building of a network of industry contacts. In addition, Lancaster University supports researchers to write for the Conversation, a news service delivering articles directly from researchers to the public; Dr Gabilon will make use of this support to produce expository articles explaining the benefits that adaptive data science approaches can deliver to society. Finally, to ensure successful public engagement, Dr Gabillon will attend Lancaster University's "The Engaging Researcher Course", a one-day experiential training course to explore public engagement activities that researchers can get involved in.

The excellent and innovative research generated in this project will of course be published Open Access in the world's leading academic journals and conferences, and all code generated will be also be released under standard Open frameworks. Prof. Leslie currently works with several companies, both large and small, including the Defence Science and Technology Laboratory who have a current interest in security games. Dr Gabillon will be mentored to develop similar relationships. He will also work with Security Lancaster (www.lancs.ac.uk/security-lancaster) to ensure the results of the current project are shared with relevant industrial and government partners. We will discuss results directly with companies in Lancaster University's Knowledge Business Centre, an innovation hub providing a gateway for business/academic interaction which allows the transfer of expertise between Lancaster's academics, regional businesses and community partnerships through training and technology transfer activities. A particularly successful mechanism deployed extensively at Lancaster is the industrially-sponsored MSc or PhD project, which allows the supervisor's research to be both developed and deployed directly within a company; Dr Gabillon will be encouraged to join appropriate supervisory teams to help both disseminate the project's research and develop an industrial research network to enhance his future career. The Research Support Office of Lancaster University has extensive experience of industrial engagement and will assist in the management of IP and any patents that may arise from the research.

## 3 Implementation

### 3.1 Overall coherence and effectiveness of the work plan, including appropriateness of the allocation of tasks and resources

**Work packages**

**WP1: Combinatorial bandits** Dr Gabillon will develop new approaches to combinatorial bandits (Objective 1), investigating the pure exploration problem and the online regret problem. This WP builds upon current research of Dr Gabillon and can be completed in months 1–6. This will result in **Deliverable 1.1**, a paper on pure exploration in combinatorial bandits, and **Deliverable 1.2**, a paper on regret in combinatorial bandits with submodular reward structures.

**WP2: Security games** Objectives 2 and 3 will be considered in this Work Package, which will develop algorithms for both pure exploration and online performance guarantees in security games. This WP brings together the knowledge of the Researcher and the Supervisor, and will therefore also be started in Month 1. However it will take longer to complete due to the greater level of novelty for Dr Gabillon, so will continue for 12 months. **Deliverable 2.1** is a paper on pure exploration in security games; **Deliverable 2.2** is a paper on online regret bounds in security games.

**WP3: Combinatorial security games** Dr Gabillon will address Objective 4 with the development of methods for security games with combinatorial action spaces. This package combines the results of WP1 and WP2. After completion of WP1, Dr Gabillon will switch attention to working in combinatorial games in parallel with WP2, and this WP will continue until the end of the project. **Deliverable 3.1** is a paper presenting the results of this research.

**WP4: Network defence** Objective 5 will be addressed, with the application of combinatorial work (WP1 and WP3) to network problems. Dr Gabillon will visit Dr Michal Valko learn about techniques in networks and how to integrate them with the previously-developed combinatorial results. He will also collaborate with researchers from Security Lancaster to develop applications in supply chain protection. This package will be started after WP2 has been completed in month 12, and run until the end of the project. **Deliverable 4.1** is a paper describing a bandit approach to network defence, and **Deliverable 4.2** is a paper describing the game-theoretical results.

**Major milestones**

**Milestone 1** is the completion of WP1. If strong results are obtained in this first phase of research, then extending to combinatorial games under an assumption that the attacker always plays a best response to the current mixed strategy will be relatively straightforward, and greater emphasis can be placed on WP3 straight away. If the results here are not so strong, more effort will need to be given to WP2 in order to obtain suitable building blocks for WP3.

**Milestone 2** is at the end of 1 year of the project, when WP1 and WP2 will both be completed, and WP3 is in progress. This will give an opportunity to take stock and decide the problems to be addressed in WP4. If the game-theoretical results in WP2 are strong, and WP3 is progressing well, then the full game-theoretical approach can be addressed directly in WP4. However weaker game-theoretical results may necessitate initial focus on bandit approaches in WP4.

### 3.2 Appropriateness of the management structure and procedures, including quality management and risk management

The Research Support Office at Lancaster University has extensive experience of managing European project grants, and will be responsible for administering the project budget, legal aspects and potential commercial exploitation of the research. Dr Gabillon will be a member of the Department of Mathematics and Statistics, and more specifically the Statistical Learning group lead by Prof. Leslie. He will also be assigned a formal mentor under standard Lancaster University human resources procedures, who will be

a second point of contact. During the project, Dr Gabillon will be responsible for the research work, and will meet weekly with Prof. Leslie to discuss results, challenges and research strategies. Dr Gabillon will also be responsible for the management of the project; he will be supervised in this task through monthly management and mentoring meetings with Prof. Leslie, in which progress against the workplan and career development plan will be discussed.
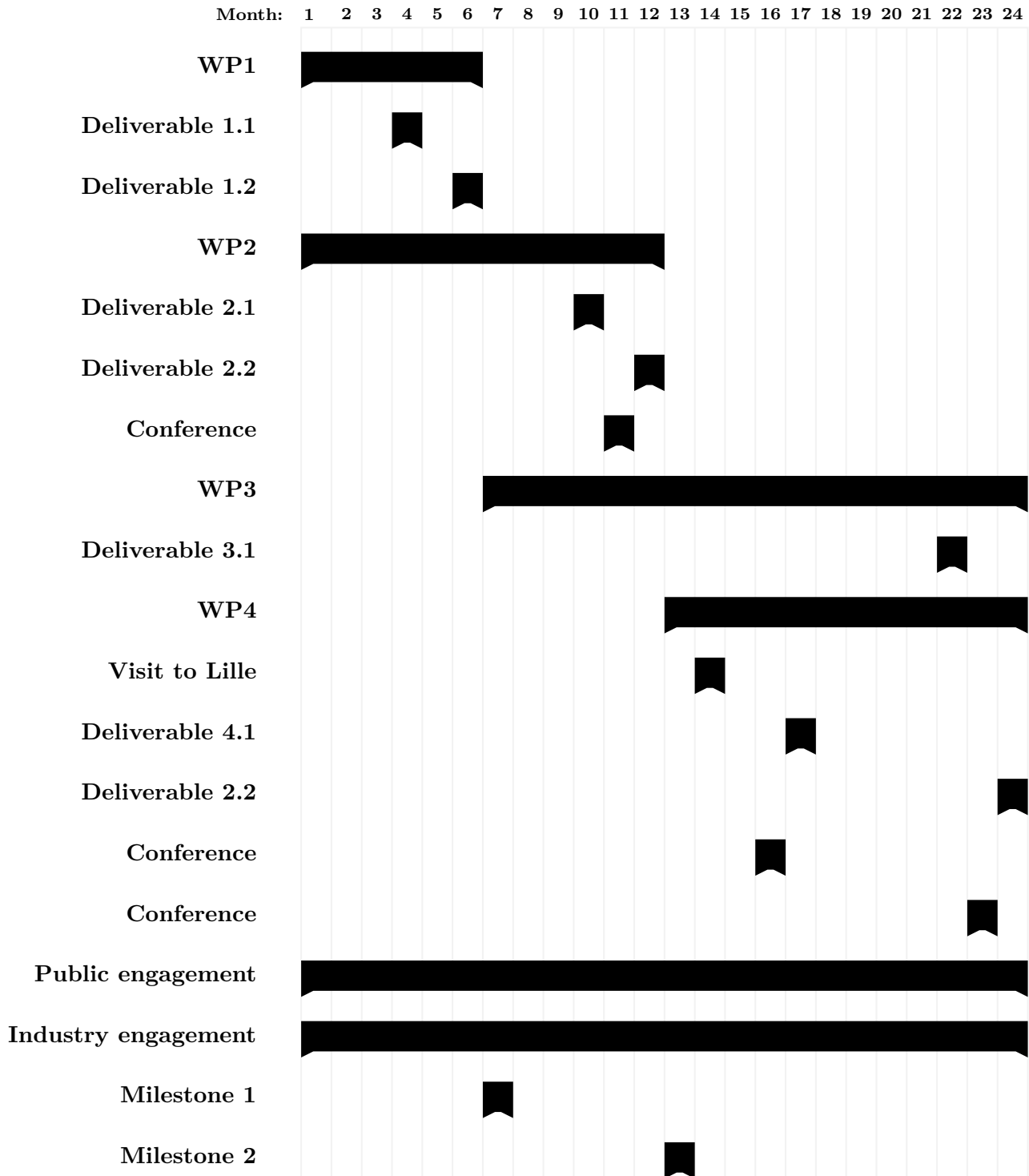
Clearly there are risks at each stage of an ambitious research project such as this. The two-pronged approach mitigates some of this risk: if developments in combinatorial approaches prove to be difficult then greater focus will be placed on the game theory, and vice versa. That being said, both of WP1 and WP2 contain elements which are lower-risk while still likely to yield high-quality research outputs. A solid foundation can thus be laid while the Researcher and Supervisor develop a working relationship, in preparation for the more ambitious objectives in the latter part of the project.

### 3.3   Appropriateness of the institutional environment (infrastructure)

Dr Gabillon will be hosted in the Department of Mathematics and Statistics, Lancaster University. Prof. Leslie will provide the main mentorship and research supervision. The Statistical Learning group, and the Statistics Research Group beyond that, will provide further immediate support to Dr Gabillon. The Department has extremely strong links with research groups in Operations Research in Lancaster University Management School, through the STOR-i Centre for Doctoral Training, and with Computer Science, through the Data Science Institute. Therefore multiple researchers in cognate areas will contribute to the project with informal mentorship and research leadership, as well as providing an environment with multiple relevant research seminars. In terms of physical resources, the Department will provide high quality office space and standard IT facilities, including high performance computing, to allow Dr Gabillon to carry out the project.

### 3.4   Competences, experience and complementarity of the participating organisations and institutional commitment

The Department of Mathematics and Statistics at Lancaster University was ranked fifth equal in the United Kingdom in the most recent Research Excellence Framework assessment. The Department has a thriving research environment, with 50 faculty, 11 post-doctoral fellows, and 72 PhD students. The Department has numerous government- and industry-funded research projects, many of which relate to industrially-motivated statistics and operations research and are related to the currently-proposed project. The skill set of Dr Gabillon complements that of the Beneficiary by providing expertise in current algorithmic approaches to bandit algorithms and reinforcement learning. The host institution in return provides expertise in statistical methodology appropriate to online inference, and game theoretical learning, and a strong track-record of working with industry to ensure the fundamental research is relevant and generates impact. In addition Dr Gabillon will develop links with Security Lancaster (www.lancaster.ac.uk/security-lancaster/), in which researchers are currently addressing the security of supply chains using game-theoretical approaches, to both develop test cases for the current research project and build links with their network of industry and government collaborators.

| Month: | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **WP1** | ████████████ | | | | | | | | | | | | | | | | | | | | | | | |
| **Deliverable 1.1** | | | ■ | | | | | | | | | | | | | | | | | | | | | |
| **Deliverable 1.2** | | | | | ■ | | | | | | | | | | | | | | | | | | | |
| **WP2** | ████████████████████████ | | | | | | | | | | | | | | | | | | | | | | | |
| **Deliverable 2.1** | | | | | | | | | ■ | | | | | | | | | | | | | | | |
| **Deliverable 2.2** | | | | | | | | | | | ■ | | | | | | | | | | | | | |
| **Conference** | | | | | | | | | | ■ | | | | | | | | | | | | | | |
| **WP3** | | | | | | ██████████████████████████████████ | | | | | | | | | | | | | | | | | | |
| **Deliverable 3.1** | | | | | | | | | | | | | | | | | | | | | | ■ | | |
| **WP4** | | | | | | | | | | | | ██████████████████████ | | | | | | | | | | | | |
| **Visit to Lille** | | | | | | | | | | | | | ■ | | | | | | | | | | | |
| **Deliverable 4.1** | | | | | | | | | | | | | | | | ■ | | | | | | | | |
| **Deliverable 2.2** | | | | | | | | | | | | | | | | | | | | | | | ■ | |
| **Conference** | | | | | | | | | | | | | | ■ | | | | | | | | | | |
| **Conference** | | | | | | | | | | | | | | | | | | | | | | ■ | | |
| **Public engagement** | ████████████████████████████████████████████████ | | | | | | | | | | | | | | | | | | | | | | | |
| **Industry engagement** | ████████████████████████████████████████████████ | | | | | | | | | | | | | | | | | | | | | | | |
| **Milestone 1** | | | | | | ■ | | | | | | | | | | | | | | | | | | |
| **Milestone 2** | | | | | | | | | | | | ■ | | | | | | | | | | | | |

# 4 CV of the Experienced Researcher

During the course of my studies several invaluable experiences have greatly contributed to my desire to pursue a research-based career in mathematics applied to machine learning. From my early years as an undergraduate student I have tried to keep the balance between theory and application. After three years of intensive Mathematics and Physics studies I entered TELECOM SudParis, a Telecommunication engineering school. There, on the one hand my engineering education made me comfortable with programming (C/C++, Java) and Network issues (LANs, WANs) and on the other and I personally got involved in a research project on PCA algorithms which has lead to a publication at ICASSP 2009. In 2008, I continued with my graduate studies in Applied Mathematics as a master student with focus on Statistical Learning where I developed solid background Machine Learning theory (including a course on Graphical Models by Francis Bach and one on Reinforcement Learning by Rémi Munos). Still I completed my master with an internship at INRIA research lab where I applied statistical learning techniques to help design a realistic automatic ad-server for Orange Inc affiliated websites. This work has launched a collaboration which is still in progress.

My current research involves the investigation of machine learning techniques to create algorithms that, in some way, adapts to its users, or more generally learns from its environment. The approach is both theoretical and application oriented. A major objective in our algorithms development is to ensure our algorithms capture the real complexity of a problem and testing in practice their performances in real world problems. During my PhD, I investigated Reinforcement Learning (RL) which is a field where one tries to solve complex systems where an agent has to learn from its environment. More precisely, the focus was on a class of algorithms called "Classification-based Policy Iteration" (CBPI) which are algorithms that learn directly the policies as output of a classifier. Thus they avoid, as in the standard RL techniques, to define a policy through an associated value function as this value function is often poorly approximated. Therefore, this class of algorithms is expected to perform better than its value-based counterparts whenever the policies are easier to represent than their value functions. However, CBPI algorithms can require large number of samples from the environment. To improve the CBPI efficiency, I proposed new hybrid approaches using value function approximations in the CBPI framework that leverage the benefits of both approaches (which led to two publications in ICML 2011 & 2012 while a journal paper has been published in JMLR). Moreover, we applied our techniques in the game of Tetris, a domain where RL techniques had obtained poor results, and learned a controller removing on average 50.000.000 lines (the best in the literature, to the best of our knowledge which is reported in a paper in NIPS 2013).

I also investigated Bandit problems. Bandit problems are the core mathematical formulation for modelling of adaptive and sequential decision-making. We designed a sampling strategy to solve several bandit problems in parallel (which led to two publications in NIPS 2011 & 2012).

During the course of my Ph.D. I worked as an research intern for 6 months at Technicolor Labs in Palo Alto California under the supervision of Branislav Kveton. Our primary goal was to improve the questionnaire asked to elicit movie preferences of users for a recommendation website. The problem was cast as an adaptive submodular maximization problem. The novelty was that we consider this problem in the case where the preferences of the users are not supposed to be known to build the questionnaire but need to be learned (which led to a publication in NIPS 2013 & AAAI 2014).

As a post-doctorate in the Queensland University of Technology, under the supervision of Peter Bartlett, I am conducting research in online learning. My first project deals with combinatorial pure exploration bandits, is set in a stochastic setting and could model network routing problem (online shortest-path problem). The second one is set in the non-stochastic (adversarial) bandit setting where the goal is to give a simple formulation of this bandit game that admits an exact minimax solution. This therefore is a more theoretical question that draws connection with game theory.

Through the experiences already described I developed my ability to work in a team environment. The international conferences, internships and summer schools I have been attending gave me the opportunity to learn and exchange with researchers from diverse horizons. In addition, teaching computer science (Algorithmic with Python & Databases) for Master and Licence students keeps enriching my communication

skills. I build up my programming skills through my curriculum in a telecommunication engineering school and later through the lectures and practical sections I gave. Moreover most of my projects have involved programming part which have made me comfortable with coding in Python and C++.

### Curriculum Vitae of the Applicant, Victor Gabillon

### Education

**PhD in Computer Science** (Accessit Award of the AI French Association, AFIA) **June 2014**
Team SequeL, INRIA Lille - Nord Europe, France
*Title:* "Budgeted Classification-based Policy Iteration"
*Domains:* Reinforcement learning & Bandits games
*Supervisors:* Mohammad Ghavamzadeh & Philippe Preux
*Examiners:*  Peter Auer (Leoben University), Olivier Cappé (Télécom ParisTech), Shie Mannor (Technion) and Csaba Szepesvári (Alberta University)

**M.Sc. Image Processing & Statistical Learning** with honours **Sep 2009**
École Normale Supérieure, Cachan, France

**Engineering Degree in Information Technology** **Sep. 2009**
TELECOM SudParis, Évry, France

### Professional Activities

**Postdoctoral Research Fellow in Statistics** *full time* **Nov 2015 − ongoing**
School of Mathematical Sciences, Queensland University of Technology, Brisbane, Australia
**PhD Researcher** *full time* **Oct 2009 − June 2014**
Team SequeL, INRIA Lille - Nord Europe, France
**Research Engineer** *full time* **Mar 2013 − Sep 2013**
Technicolor Research Group, Palo Alto, USA.
**External Lecturer** *part time* **Fall 2012**
Lille 1 University, France
**External Lecturer** *part time* **2010 − 2011**
Lille 3 University, France
**Research Engineer** *full time* **June 2008 − Sep 2008**
Chinese Academy of Science, Beijing, China.

### Awards & Grants

**Postdoctoral Research Fellowship in Statistics** **Nov 2015**
Two-year fellowship funded by the Queensland University of Technology
**Second place award for the best French PhD in Artificial Intelligence** **June 2015**
Award from AFIA, the French Association for Artificial Intelligence.
**Best applied paper award** **Jan 2010**
Award from the EGC conference, French speaking conference on knowledge mining and management.
**PhD Grant** **Oct 2009**
Three-year grant funded by the French Ministry of Research

### Research Expeditions

**3 months at Berkeley Statistic Departement, USA** **Mar − June 2015**
Hosted by Peter Bartlett
**One week at Inria Nancy-Grand Est, France** **June 2012**
Hosted by Bruno Scherrer of the team Maia

## Peer Reviewer

I have been an official reviewer for the Neural Information Processing Systems (NIPS) international conference in 2014 and 2015 and I have reviewed papers for the Machine Learning Jounal and the Journal of Machine Leaning Research (JMLR).

## Invited Presentations

### Talks other than Conference presentations
**Talk** Oxford Robotics Research Group Seminar, Oxford, UK, May 2014
"Classification-Based Policy Iteration perform well in the game of Tetris".
**Talk** Gatsby Reinforcement Learning Research Group, London, UK, May 2014
"Classification-Based Policy Iteration perform well in the game of Tetris".
**Talk** Team Maia Seminar, Nancy, France, June 2012
"Pure Exploration Bandits".
**Talk** Co-Adapt Seminars, Marseille, France, May 2012
"Pure Exploration Bandits for Brain-Computer Interface?".

## Publications

### Peer-reviewed journal article
J1.  Bruno Scherrer, Mohammad Ghavamzadeh, Victor Gabillon & Matthieu Geist, ***Approximate Modified Policy Iteration***, to appear in Journal of Machine Learning Research (JMLR).

### Peer-reviewed conference articles
C9.  Victor Gabillon, Branislav Kveton, Zheng Wen, Brian Eriksson & S. Muthukrishnan, ***Large Scale Optimistic Adaptive Submodularity***. AAAI 2014, $28^{th}$ Conference of the Association for the Advancement of Artificial Intelligence. Oral presentation at Quebec City, Canada, July 2014.

C8.  Victor Gabillon, Mohammad Ghavamzadeh & Bruno Scherrer, ***Approximate Dynamic Programming Finally Performs Well in the Game of Tetris***. NIPS 2013, $27^{th}$ Conference on Neural Information Processing Systems. Poster presentation at South Lake Tahoe, Nevada, December 2013.

C7.  Victor Gabillon, Branislav Kveton, Zheng Wen, Brian Eriksson & S. Muthukrishnan, ***Adaptive Submodular Maximization in Bandit Setting***. NIPS 2013, $27^{th}$ Conference on Neural Information Processing Systems. Poster presentation at South Lake Tahoe, Nevada, December 2013.

C6.  Victor Gabillon, Mohammad Ghavamzadeh & Alessandro Lazaric, ***Best Arm Identification: A unified approch to fixed budget and fixed confidence***. NIPS 2012, $26^{th}$ Conference on Neural Information Processing Systems. Poster presentation at South Lake Tahoe, Nevada, December 2012.

C5.  Bruno Scherrer, Mohammad Ghavamzadeh, Victor Gabillon & Matthieu Geist, ***Approximate Modified Policy Iteration***. ICML 2012, $29^{th}$ International Conference on Machine Learning. Long lecture presentation at Edinburgh, Scotland, June 2012.

C4.  Victor Gabillon, Mohammad Ghavamzadeh, Alessandro Lazaric & Sébastien Bubeck, ***Multi-Bandit Best Arm Identification***. NIPS 2011, $25^{th}$ Conference on Neural Information Processing Systems. Poster presentation at Granada, Spain, December 2011.

C3.  Victor Gabillon, Alessandro Lazaric, Mohammad Ghavamzadeh & Bruno Scherrer, ***Classification-based Policy Iteration with a Critic***. ICML 2011, $28^{th}$ International Conference on Machine Learning. Lecture presentation at Bellevue, USA, June 2011.

C2.  Victor Gabillon, Jérémie Mary & Philippe Preux, ***Affichage de publicités sur des portails web***. EGC 2010, $10^{th}$ French-speaking International Conference on Knowledge Extraction and Management. Lecture presentation of long article at Hammamet, Tunisia, January 2010. Best applied paper award.

C1.  Jean-Pierre Delmas & Victor Gabillon, ***Asymptotic performance analysis of PCA algorithms based on the weighted subspace criterion***. ICASSP 2009, International Conference on Acoustics, Speech and Signal Processing. Poster presentation at Taipei, Taiwan, April 2009.

*Peer-reviewed workshop article*

W1.   Victor Gabillon, Alessandro Lazaric & Mohammad Ghavamzadeh, ***Rollout Allocation Strategies for Classification-based Policy Iteration***. Workshop on Reinforcement Learning and Search in Very Large Spaces International Conference on Machine Learning, Lecture presentation at Haifa, Israel, June 2010.

*Major research achievements & industrial innovations*

- *Research: **Reinforcement Learning is Finally Competitive:*** We proposed a new family of reinforcement learning methods based on "Classification-based Policy Iteration" algorithms. In addition to providing thorough theoretical analysis of these methods (C3,C5), we implemented the proposed algorithms and ran extensive experimental studies to analyse their performance using the game of Tetris as a benchmark. Our results show an unprecedented performance of reinforcement learning methods in Tetris, improving upon the state-of-the-art techniques. Moreover, while these state-of-the-art techniques were based on black-box optimisation methods that require a large number of samples from the environment, our methods are able to learn Tetris strategies and achieve the same performance with 10 times less number of samples (C8).

- *Industry: **Constrained Learning for Orange Ad Server:*** Orange, the leading company to provide telecommunication solutions in France, had made a contract with the research team SequeL in order to automate their online web-advertising services. My initial goal was to make a survey of the machine learning literature and find an appropriate solution to optimise their click-through rate revenues. This solution had to take into account specific new constraints on the limited and known number of display per ads. I proposed a new approach combining linear programming and bandit algorithms and validated its performance on synthetic data. The results received the best paper award at the French conference on Data Extraction and Knowledge Management (C2). Moreover, the project initiated an ongoing collaboration between SequeL and Orange.

- *Industry: **Adaptive Questionnaire Design at Technicolor Inc:*** During the course of my PhD, I participated in a 6-months R&D internship program at Technicolor Inc.'s research lab. The primary goal of my project was to improve upon an online questionnaire which aimed to elicit users' movie preferences and generate a recommender system. I cast the program as an adaptive submodular maximisation problem. The novelty of my approach was to consider a completely realistic formulation of the problem where the users' preferences were unknown and to be actively learned in order to build the adaptive questionnaire. My efforts in this short period of time resulted in the publication of two peer-reviewed papers at prestigious international conferences in machine learning (C7, C9).

**Teaching**

In the past 5 years I taught 216 hours of undergraduate and master's courses in France.

*Instructor:*

- *Introduction to algorithmic and programming with Python.*
  48 hours (lectures and practical sessions). Winter 2010, Fall 2011 & Fall 2012
  $1^{rst}$ year of Master *Computer science and document* at Lille 3 University and $1^{rst}$ year of Licence *Physics-Chemistry* at Lille 1 University.

*Teaching assistant:*

- *SQL and Python.* 36 hours (practical sessions). Fall 2010.
  $3^{rd}$ year of Licence *Mathematics and computer science applied to social sciences* at Lille 3 University.

- *Designing databases and object-oriented programming.* 36 hours (practical sessions). Winter 2011.
  $3^{rd}$ year of Licence *Mathematics and computer science applied to social sciences* at Lille 3 University.

# 5   Capacities of the Participating Organisations

**Beneficiary: Lancaster University**

| | |
|---|---|
| **General Description** | Lancaster University is a top ten UK university. The Department of Mathematics and Statistics, within the Faculty of Science and Technology, hosts one of the largest and strongest statistics research groups in the UK comprising 25 academic staff, 10 research associates and around 50 FTE research students. In the 2014 Research Excellence Framework assessment, the Mathematical Sciences at Lancaster were ranked fifth overall and third in terms of the impact of research. Research is supported by grants from the UK Research Councils, the European Commission, and industrial sponsors. The statistics research group is also a fundamental partner in Lancaster's new Data Science Institute, which aims to act as a catalyst for Data Science, providing an end-to-end interdisciplinary research capability. |
| **Role and Commitment of key persons (supervisor)** | Prof. David Leslie, PhD in Mathematics (University of Bristol, 2003). 17 PhD students and 2 post-doctoral fellows supervised. 5% FTE time commitment to the project throughout the 24 month duration. |
| **Key Research Facilities, Infrastructure and Equipment** | The Department of Mathematics and Statistics is housed in dedicated space at Lancaster University. The Researcher will be provided with office space and basic equipment within the Department. Researchers in have access to the Department's own computer support (2.6FTE computer technicians) and computer cluster (nearly 500 computer cores, 800GB of memory). These computing facilities are supplemented by access to Lancaster University's High-End Computing cluster (1700 computer cores, 8TB of memory, 32TB of high performance filestore). |
| **Independent research premises?** | Yes |
| **Previous Involvement in Research and Training Programmes** | Between 2001 and 2005 the department held the Marie Curie Training Site status for its PhD programme. The Postgraduate Statistics Center (PSC) was founded in 2005 as the only Centre for Excellence in Teaching and Learning focussing on postgraduate statistics in the UK. The PSC is still operative and runs three Masters degrees (Statistics, Quantitative Methods, and Quantitative Finance) and coordinates the PhD programme in statistics. |
| **Current involvement in Research and Training Programmes** | Together with the Management School, the Department hosts and runs STOR-i, a multimillion pound EPSRC-funded Centre for Doctoral Training in Statistics and Operational Research in partnership with industry. The Centre was established in 2010 and funds 12 PhD students per year. The department is also a key player in the Academy for Phd Training in Statistics, a collaboration between major UK statistics research groups to organise courses for first-year PhD students in statistics and applied probability nationally. The group hosts one node of a multi-institution Programme Grant on Intractable Likelihood, and received industrial funding from companies including Shell, BT, Google and Unilever. The Department's Medical and Pharmaceutical Statistics Research Unit works closely with the pharmaceutical industry and public sector research institutes to develop novel statistical methods for the design and analysis of clinical trials. It leads the EU-funded research training network IDEAS (www.ideas-itn.eu) and is an integral part of the Medical Research Council funded North-West Hub for Trials Methodology Research. |
| **Relevant Publications and/or research/innovation products** | Perkins, S. and Leslie, D.S. (2014) Stochastic fictitious play with continuous action sets. *Journal of Economic Theory* **152**, 179–213.<br>Chapman, A.C., Leslie, D.S., Rogers, A. and Jennings, N.R. (2013) Convergent learning algorithms for unknown reward games. *SIAM Journal on Control and Optimization* **51**, 3154-3180.<br>May, B.C., Korda, N., Lee, A. and Leslie, D.S. (2012) Optimistic Bayesian sampling in contextual-bandit problems. *Journal of Machine Learning Research* **13**, 2069–2106.<br>Larsen, T., Leslie, D.S., Collins, E.J. and Bogacz, R. (2010) Posterior weighted reinforcement learning with state uncertainty. *Neural Computation* **22**, 1149–1179.<br>Leslie, D.S. and Collins, E.J. (2003) Convergent multiple-timescales reinforcement learning algorithms in normal form games. *Annals of Applied Probability* **13**, 1231–1251. |

## ENDPAGE


MARIE SKLODOWSKA-CURIE ACTIONS


## Individual Fellowships (IF)
## Call: H2020-MSCA-IF-2014


PART B


"OSEGA"


## This proposal is to be evaluated as:

## [Standard EF]