A THEORY OF MULTI-AGENT GENERATIVE FLOW NETWORKS

Anonymous authors

Paper under double-blind review

Abstract

Generative flow networks utilize a flow-matching loss to learn a stochastic policy for generating objects from a sequence of actions, such that the probability of generating a pattern can be proportional to the corresponding given reward. However, a theoretical framework for multi-agent generative flow networks (MA-GFlowNets) has not yet been proposed. In this paper, we propose the theory framework of MA-GFlowNets, which can be applied to multiple agents to generate objects collaboratively through a series of joint actions. We further propose four algorithms: a centralized flow network for centralized training of MA-GFlowNets, an independent flow network for decentralized execution, a joint flow network for achieving centralized training with decentralized execution, and its updated conditional version. Joint Flow training is based on a local-global principle allowing to train a collection of (local) GFN as a unique (global) GFN. This principle provides a loss of reasonable complexity and allows to leverage usual results on GFN to provide theoretical guarantees that the independent policies generate samples with probability proportional to the reward function. Experimental results demonstrate the superiority of the proposed framework compared to reinforcement learning and MCMC-based methods.

026 027 028

029

025

004

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

Generative flow networks (GFlowNets) Bengio et al. (2023) can sample a diverse set of candidates in an active learning setting, where the training objective is to approximate sampling of the candidates proportionally to a given reward function. Compared to reinforcement learning (RL), where the learned policy is more inclined to sample action sequences with higher rewards, GFlowNets can perform exploration tasks better. The goal of GFlowNets is not to generate a single highest-reward action sequence, but rather is to sample a sequence of actions from the leading modes of the reward function Bengio et al. (2021). However, based on current theoretical results, GFlowNets cannot support multi-agent systems.

A multi-agent system is a set of autonomous interacting entities that share a typical environment, perceive through sensors, and act in conjunction with actuators Busoniu et al. (2008). Multi-agent 040 reinforcement learning (MARL), especially cooperative MARL, is widely used in robot teams, dis-041 tributed control, resource management, data mining, etc Zhang et al. (2021); Canese et al. (2021); 042 Feriani & Hossain (2021). There two major challenges for cooperative MARL: scalability and 043 partial observability Yang et al. (2019); Spaan (2012). Since the joint state-action space grows ex-044 ponentially with the number of agents, coupled with the environment's partial observability and communication constraints, each agent needs to make individual decisions based on the local action observation history with guaranteed performance Sunehag et al. (2018); Wang et al. (2020); 046 Rashid et al. (2018). In MARL, to address these challenges, a popular centralized training with 047 decentralized execution (CTDE) paradigm Oliehoek et al. (2008); Oliehoek & Amato (2016) is pro-048 posed, in which the agent's policy is trained in a centralized manner by accessing global information and executed in a decentralized manner based only on the local history. However, extending these techniques to GFlowNets is not straightforward, especially in constructing CTDE-architecture flow 051 networks and finding IGM conditions for flow networks need investigating. 052

1053 In this paper, we propose the multi-agent generative flow networks (MA-GFlowNets) framework for cooperative decision-making tasks. Our framework can generate more diverse patterns through

054 sequential joint actions with probabilities proportional to the reward function. Unlike vanilla GFlowNets, the proposed method analyzes the interaction of multiple agent actions and shows how 056 to sample actions from multi-flow functions. Our approach consists of building a virtual global GFN 057 capturing the policies of all agents and ensuring consistency of local (agent) policies. Variations of 058 this approach yield different flow-matching losses and training algorithms.

Furthermore, we propose the Centralized Flow Network (CFN), Independent Flow Network (IFN), 060 Joint Flow Network (JFN), and Conditioned Joint Flow Network (CJFN) algorithms for multi-agent 061 GFlowNets framework. CFN considers multi-agent dynamics as a whole for policy optimization 062 regardless of the combinatorial complexity and demand for independent execution, so it is slower; 063 while IFN is faster, but suffers from the flow non-stationary problem. In contrast, JFN and CJFN, which are trained based on the local-global principle, takes full advantage of CFN and IFN. They can 064 reduce the complexity of flow estimation and support decentralized execution, which are beneficial 065 to solving practical cooperative decision-making problems. 066

067 Main Contributions: 1) We first generalize the measure GFlowNets framework to the multi-agent 068 setting, and propose a theory of multi-agent generative flow networks for cooperative decision-069 making tasks; 2) We propose four algorithms under the measure framework, namely CFN, IFN, JFN and CJFN, for training multi-agent GFlowNets, which are respectively based on centralized 071 training, independent execution, and the latter two algorithms are based on the CTDE paradigm; 3) We propose a local-global principle and then prove that the joint state-action flow function can be 072 decomposed into the product form of multiple independent flows, and that a unique Markovian flow 073 can be trained based on the flow matching condition; 4) We conduct experiments based on cooper-074 ative control tasks to demonstrate that the proposed algorithms can outperform current cooperative 075 MARL algorithms, especially in terms of exploration capabilities. 076

077 078

079

2 **PROBLEM FORMULATION**

The **multi-agent setting** formalizes the data of state, actions and transitions for multiple agents. 080 The state space S as well as the state-dependent action spaces A_s are measurable spaces; for 081 each state $s \in S$, the environment comes with a stochastic transition map¹ $\mathcal{A}_s \xrightarrow{T_s} S$. We formalize this dependency on state by bundling (packing) state and action together into a bundle 082 083 $\{(s,a) \mid s \in S, a \in A_s\} = A \xrightarrow{S,T} S$ where $S(s,a) \coloneqq s$ and $T(s,a) \coloneqq T_s(a)$. For graphs, a bundled action is an edge $s \to s'$, the statemap S returns the origin s while the transition map returns 084 085

086 the destination s'. A policy is then a section of S is a kernel $\mathcal{S} \xrightarrow{\pi} \mathcal{A}$ such that $S \circ \pi$ is identity on \mathcal{S} . 087 Each agent $i \in I$ in the finite agent set I has its own observation 088 $o^{(i)}$ in its observation space $\mathcal{O}^{(i)}$; it depends on the state via the 089 projection $\mathcal{S} \xrightarrow{p^{(i)}} \mathcal{O}^{(i)}$. For simplicity sake, we identify \mathcal{S} = 090 $\prod_{i \in I} \mathcal{O}^{(i)}$. Each agent has its own action space $\mathcal{A}^{(i)}$ and each 091 of the agent observation-dependent action space A_o contains a 092 special action STOP; the environment is such that once an agent 093 chooses STOP, it is put on hold until all agents do as well. The 094 game finishes when all agent have chosen STOP, a reward is 095 given based on the last state. The reward received is formalized 096



Figure 1: Multi-agent formalism

by a non-negative function $r: S \to S$. We assume that each agent may freely choose its own action independently from the actions chosen by other agents: this is formalized via $A_s = \prod_{i \in I} A_{o^{(i)}}^{(i)} / \sim$ ie the Cartesian product of agent actions space up to identification of the STOP actions. A trajectory 098 099 of the system of agents is a, possibly infinite, sequence of states $(s_t)_{t < \tau+1}$ with $\tau \in \mathbb{N} \cup \{\infty\}$ starting 100 at the source state $s_0 \in S$ and may eventually calling STOP; the space of trajectories is \mathcal{T} . A policy 101 on S induces a Markov chain hence a distribution on trajectories. 102

Our objective is to build a policy π *so that the induced trajectories are finite and* s_{τ} *is distributed* 103 proportionally to $R \coloneqq r\lambda$ where λ is some fixed measure on S and $\int_{s \in S} r(s) d\lambda(s)$ is finite. 104

¹⁰⁵ ¹We adopt the naming convention of Douc et al. (2018). The kernel $K: \mathcal{X} \to \mathcal{Y}$ is a stochastic map which 106 is formalized as follows: for all $x \in \mathcal{X}, K(x \to \cdot)$ is a probability distribution on \mathcal{Y} . In addition, $K(x \to \cdot)$ varies measurably with x in the sense that for all measurable set $A \subset \mathcal{Y}$, the real valued map $x \mapsto K(x \to A)$ 107 is measurable.

Measurable GFlowNets Brunswic et al. (2024); Lahlou et al. (2023); Li et al. (2023d); Deleu & Bengio (2023); Bengio et al. (2023) are defined in the single-agent setting i.e.

- 111
- 112 113

114

122 123 124

137 138 139

150 151 152

153

154 155 156 $\mathcal{A} \underbrace{\overbrace{\tau}}^{S}_{T} \mathcal{S} \xrightarrow{R} \mathbb{R}_{+}, \text{ with } |I| = 1.$

115 A GFlowNets on (S, A, S, T, R) is a forward policy $\pi : S \to A$ together with a non-negative finite 116 measure F_{out} on S called the outflow or state-flow. The reward is generally non-trainable and 117 unknown but implicitly a component of F_{out} and π ; since the reward may not be tractable in the 118 multi-agent setting, we favor a reward-free parameterization of GFlowNets. We thus parameterize 119 them by triplets $(\pi^*, F_{out}^*, F_{init})$ where $\pi^*(s) = \pi(s \mid a \neq \text{STOP}), F_{out}^* := F_{out} - R$ and $F_{init} =$ 120 $F_{out}(s_0)\pi(s_0)$. The *-notations informally mean we restrict the objects to $S \setminus \{s_0, s_f\}$. Given a GFlowNet in *-notations together with a reward, there is a unique GFlowNet in usual notations.

A GFlowNet is trained to satisfy the so-called flow-matching constraint:

$$F_{\text{out}} = F_{\text{in}} \coloneqq F_{\text{init}} + F_{\text{out}}^* \pi^* T, \tag{1}$$

as measures on S. In passing we introduce $\hat{R} \coloneqq F_{in} - F_{out}^*$, $F_{in}^* \coloneqq F_{out}^* \pi^* T$ and $F_{action} \coloneqq F_{out} \otimes \pi$. The induced Markov chain starts at s_0 sampled from the unnormalized distribution F_{init} and then for every t the policy is applied until the action STOP is picked: $a_t \sim \pi(s_t \to \cdot)$ and if $a_t \neq STOP$, $s_{t+1} \sim T(a_t \to \cdot)$. Usually we choose $F_{init} \propto \ell(\theta)$ with ℓ a known, easily sampled from, distribution family. The sampling time τ is the t such that $a_t = STOP$.

The following Theorem was first proved on DAG in Bengio et al. (2023) and shows GFlowNets
 answer our problem definition.

Theorem 1 ((Brunswic et al., 2024) Theorem 2) Let $\mathbb{F} := (\pi, F_{out}^*, F_{init})$ be a GFlowNets on (S, A, S, T, R). If the reward R is non-zero and \mathbb{F} satisfies the flow-matching constraint, then its sampling time is almost surely finite and the sampling distribution is proportional to R. More precisely:

$$\mathbb{P}(\tau < +\infty) = 1, \qquad \mathbb{E}(\tau) \le \frac{F_{\text{out}}(\mathcal{S})}{R(\mathcal{S})} - 1, \quad and \quad s_{\tau} \sim \frac{1}{R(\mathcal{S})}R.$$
(2)

Flow-matching losses (FM), denoted by \mathcal{L}_{FM} , compare the outflow F_{out} with the inflow $F_{in} := F_{init} + F_{out}T\pi$; They are minimized when $F_{in} = F_{out}$ so that, surely, a gradient descent on GFlowNets parameters enforces equation 1. In the original works Bengio et al. (2021); Malkin et al. (2022), Bengio et al. used divergence-based FM losses valid as long as the state space does not have cycle and Brunswic et al. (2024) introduced stable FM losses allowing training in presence of cycles:

$$\mathcal{L}_{\rm FM}^{\rm div}(\mathbb{F}^{\theta}) = \mathbb{E}_{s \sim \nu_{\rm state}} g \circ \log\left(\frac{dF_{\rm in}^{\theta}}{dF_{\rm out}^{\theta}}(s)\right)$$
(3)

$$\mathcal{L}_{\rm FM}^{\rm stable}(\mathbb{F}^{\theta}) = \mathbb{E}_{s \sim \nu_{\rm state}} g\left(\frac{dF_{\rm in}^{\theta}}{d\lambda}(s) - \frac{dF_{\rm out}^{\theta}}{d\lambda}(s)\right),\tag{4}$$

where g is some positive function, decreasing on $[-\infty, 0]$, g(0) = 0 and increasing on $[0, +\infty]$. A practical stable training loss on graphs can be written as

$$\mathcal{L}(\mathbb{F}^{\theta}) = \mathbb{E}\sum_{t=1}^{\tau} \left\{ \log \left[1 + \varepsilon \left| F_{\text{in}}^{\theta} \left(s_{t} \right) - F_{\text{out}}^{\theta} \left(s_{t} \right) \right|^{\alpha} \right] \times \left(1 + \eta \left(F_{\text{in}}^{\theta} \left(s_{t} \right) + F_{\text{out}}^{\theta} \left(s_{t} \right) \right) \right)^{\beta} \right\},$$
(5)

where s_t are path sampled from **any** distribution of paths in S, and the parameters satisfy the condition { $\varepsilon, \eta, \alpha, \beta > 0$ }.

160 **MA-GFlowNets** are tuples $((\mathbb{F}^{(i)})_{i \in I}, \mathbb{F})$, where each *local* GFlowNets $\mathbb{F}^{(i)}$ is defined on 161 $(\mathcal{O}^{(i)}, \mathcal{A}^{(i)}, S^{(i)}, T^{(i)}, R^{(i)})$ for $i \in I$ and the *global* GFlowNets \mathbb{F} is defined on $(\mathcal{S}, \mathcal{A}, S, T, R)$. In general, some GFlowNets (local or global) may be virtual in the sense that it is not implemented.

162 MULTI-AGENT GFLOWNETS 3

163 164

This section is devoted to details and theory regarding the variations of algorithms for MA-GFlowNets training. If resources allow, the most direct approach is included in the training of 166 the global model directly, leading to a centralized training algorithm in which the local GFlowNets 167 are virtual. As expected, such an algorithm suffers from high computational complexity, hence, 168 demanding decentralized algorithms. Decentralized algorithms require the agents to collaborate to some extent. We achieve such a collaboration by enforcing consistency rules between the local and 169 170 global GFlowNets. The global GFlowNets is virtual and is used to build a training loss for the local models ensuring the global model is GFlowNets, so that the sampling Theorem applies. The sam-171 pling properties of the MA-GFlowNets are then deduced from the flow-matching property of the 172 virtual global model.

173 174 175

176

177

178

179

181 182 183

192

193

194

195

196 197

198 199

201

3.1 CENTRALIZED TRAINING

Centralized training consists in training of the global flow directly. Here, the local flows are virtual in the sense that they are recovered from the global flow as image by the observation maps. We use FM-losses as given in equations 3-4 applied to the flow on (S, A). See Algorithm 1. Implicitly, $F_{\rm out}$ contains a parameterizable component from $F_{\rm out}^{*}$, while $F_{\rm in}$ contains the parameterization of π^* and F_{init} .

Algorithm 1 Centralized Flow Network Training Algorithm for MA-GFlow	wNets
Input: A multi-agent environment $(S, A, O^{(i)}, A^{(i)}, p_i, S, T, R)$, a parameter	eterized GFlowNets F :=
$(\pi, F_{\text{out}}^*, F_{\text{init}})$ on $(\mathcal{S}, \mathcal{A})$.	
while not converged do	
Sample and add trajectories $(s_t)_{t\geq 0} \in \mathcal{T}$ to replay buffer with policy	$\pi(s_t \to a_t).$
Generate training distribution ν_{state} .	
Apply minimization step of the FM loss $\mathcal{L}^{ ext{stable}}_{ ext{FM}}(\mathbb{F}^{ heta})$.	
end while	

From the algorithmic viewpoint, the CFN algorithm is identical to a single GFlowNets. As a consequence, usual results on the measurable GFlowNets apply as is. There are, however, a number of key difficulties: 1) even on graphs, the computational complexity increases as $O(|\mathcal{A}_s|^N)$ at any given explored state; 2) centralized training requires all agents to share observations, which is impractical since in many applications the agents only have access to their own observations.

LOCAL TRAINING: INDEPENDENT 3.2

The dual training method is embodied in the training of local GFlowNets instead of the global one. In this case, the local flows $\mathbb{F}^{(i)}$ are parameterized and the global flow is virtual. In the same way, 200 a local FM loss is used for each client. In order to have well-defined local GFlowNets, we need a local reward, for which a natural definition is $R^{(i)}(o_t^{(i)}) := \mathbb{E}(R(s_t)|o_t^{(i)})$. The local training loss 202 function can be written as:

$$\mathcal{L}(\mathbb{F}^{(i)}) = \mathbb{E}\sum_{t=1}^{\tau} \left\{ \log \left[1 + \varepsilon \left| F_{\text{in}}^{\theta_i} \left(o_t^i \right) - F_{\text{out}}^{\theta_i} \left(o_t^i \right) \right|^{\alpha} \right] \times \left(1 + \eta \left(F_{\text{in}}^{\theta_i} \left(o_t^i \right) + F_{\text{out}}^{\theta_i} \left(o_t^i \right) \right) \right)^{\beta} \right\}.$$
(6)

207 The algorithm 3 in Appendix B describes a simplest training 208 method, which solves the issue of exponential action com-209 plexity with an increasing number of agents. In this formulation, however, two issues arise: the evaluation of ingoing flow 210 $F_{in}^{(i)}(o^{(i)})$ becomes harder as we need to find all transitions 211 leading to a given local observation (and not to a given global 212 state). This problem may be non-trivial as it is also related 213 to the actions of other agents. More importantly, in this case, 214 the local reward is intractable so we cannot accurately estimate 215 the reward $R^{(i)}(o^{(i)})$ of each node; Falling back to using the



Figure 2: Performance comparison on Hyper-grid task.

216 stochastic reward $R^{(i)}(o^{(i)}) \coloneqq R(s_t|o_t^{(i)})$ instead leads to transition uncertainty and spurious re-217 wards, which can cause non-stationarity and/or mode collapse as shown in Figure 2. 218

3.3 LOCAL-GLOBAL TRAINING

219

220 221

222

223

224 225

226

227

232 233

234 235 236

237

238

239

240

249 250

251

252 253 254

255

256

257

258

3.3.1 LOCAL-GLOBAL PRINCIPLE: JOINT FLOW NETWORK

Local-global training is based upon the following local global principle which combined with Theorem 1 ensures that the MA-GFlowNet have sampling distribution proportional to the reward R.

Theorem 2 (Joint MA-GFlowNets) Given local GFlowNets $\mathbb{F}^{(i)}$ on some environments $(\mathcal{O}^{(i)}, \mathcal{A}^{(i)}, S^{(i)}, T^{(i)})$ there exists a global GFlowNets $\mathbb{F}^{\text{joint}}$ on a multi-agent environment $(\prod_{i \in I} \mathcal{O}^{(i)}, \mathcal{A}, S, \tilde{T})$ consistent with the local GFlowNets $\mathbb{F}^{(i)}$, such that

$$F_{\text{out}}^* = \prod_{i \in I} F_{\text{out}}^{(i),*}, \quad F_{\text{in}} = \prod_{i \in I} F_{\text{in}}^{(i)}.$$
 (7)

Moreover, if $\mathbb{F}^{\text{joint}}$ satisfies equation 1 for a reward R and each $\hat{R}^{(i)} \ge 0$ then $R = \prod_{i \in I} \hat{R}^{(i)}$.

Our Joint Flow Network (JFN) algorithm, leverage Theorem 2 by sampling trajectories with policy

$$o_t^{(i)} = p_i(s_t^{(i)}) \text{ and } \pi^{(i)}(o_t^{(i)} \to a_t^{(i)}), \ i \in I$$
 (8)

with $a_t = (a_t^{(i)} : i \in I)$ and $s_{t+1} = T(s_t, a_t)$, build formally the (global) joint GFlowNet from local GFlowNets and train the collection of agent via the FM-loss of the joint GFlowNet. Equation 7 ensures that the inflow and outflow of the (global) joint GFlowNet are both easily computable from the local inflows and outflows provided by agents. See algorithm 2.

Algorithm 2 Joint Flow Network Training Algorithm for MA-GFlowNets
Input: Number of agents N, A multi-agent environment $(S, A, O^{(i)}, A^{(i)}, p_i, S, T, R)$.
Input: Local parameterized GFlowNets $(\pi^{(i),*}, F_{out}^{(i),*}, F_{init}^{(i)})_{i \in I}$.
while not converged do
Sample and add trajectories $(s_t)_{t\geq 0} \in \mathcal{T}$ to replay buffer with policy according to equation 8.
Generate training distribution of states ν_{state} from the replay buffer.
Apply minimization step of the FM loss $\mathcal{L}_{\text{FM}}^{\text{stable}}(\mathbb{F}^{\theta,\text{joint}})$ for reward R.
and while

This training regiment presents two key advantages: over centralized training, the action complexity is linear w.r.t. the number of agents and local actions as in the independent training; over independent training, the reward is not spurious. Indeed, in $\mathcal{L}_{FM}^{\text{stable}}(\mathbb{F}^{\theta,\text{joint}})$, by equation 7, the computation of F_{in} and F_{out}^* reduces to computing the inflow and star-outflow for each local GFlowNets. Also, only the global reward R appears. The remaining, possibly difficult, challenge is the estimation of local ingoing flows from the local observations as it depends on the local transitions $T^{(i)}$, see first point below. At this stage, the relations between the global/joint/local flow-matching constraints are unclear, and furthermore, the induced policy of the local GFlowNets still depends on the yet undefined local rewards. The following point clarify those links.

259 First, the collection of local GFlowNets induces local transitions kernels $T^{(i)} : \mathcal{O}^{(i)} \to \mathcal{O}^{(i)}$ 260 which are not uniquely determined in general by a single GFlowNets. Indeed, the local policies induce a global policy $\pi(s_t \to a_t) := \prod_{i \in I} \pi(o_t^{(i)} \to a_t^{(i)})$. Then, the (virtual) transition kernel 261 262 $T^{(i)}(a_t^{(i)}) = (T(a_t)|a_t^{(i)})$ of the GFlowNets *i* depends on the distribution of states and the corre-263 sponding actions of **all** local GFlowNets. See appendix A.5 for details. Note that $T^{(i)}$ are derived 264 from the actual environment T and the joint GFlowNets on the multi-agent environment with the 265 true transition T, while the Theorem above ensures splitting of star-inflows and virtual rewards only 266 for the approximated T. Furthermore, local rewards may be formalized as a stochastic reward to 267 take into account the lack of information of a single agent, but they are never used during training: 268 the allocation of rewards across agents is irrelevant. Only the virtual rewards $\hat{R}^{(i)} = F_{out}^{(i),*} - F_{in}^{(i)}$ are relevant but they are effectively free. As a consequence, Algorithm 2 effectively trains both the 269

270 joint flow as well as a product environment model. But since in general $T \neq \tilde{T}$ Algorithm 2 may fail 271 to reach satisfactory convergence. 272

Second, beware that in our construction of the joint MA-GFlowNets, there is no guarantee that the 273 global initial flow is split as the product of the local initial flows. In fact, we favor a construction in 274 which F_{init} is non-trivial to account for the inability of local agents to assess synchronization with 275 another agent. See Appendix A.8 for formalization details. 276

Third, we may partially link local and global flow-matching properties. 277

278 **Theorem 3** Let $(\mathbb{F}^{(i)})_{i \in I}$ be local GFlowNets and let \mathbb{F} be their joint GFlowNets. Assume that none of the local GFlowNets are zero and that each $\hat{R}^{(i)} \ge 0$. If \mathbb{F} satisfies equation 1, then there exists an 280 "essential" subdomain of each $\mathcal{O}^{(i)}$ on which local GFlowNets satisfy the flow-matching constraint.

282 The restriction regarding the domain on which local GFlowNets satisfy the flow-matching constraint 283 is detailed in Appendix A.8, this sophistication arises because of the stopping condition of the multi-284 agent system. The essential domain may be informally formulated as "where the local agent is still 285 playing": an agent may decide (or be forced) to stop playing, letting other agents continue playing, 286 the forfeited player is then on hold until the game stops and rewards are actually awarded.

287 To conclude, the joint GFlowNets provides an approximation of the target global GFlowNets, this 288 approximation may fail if the transition kernel T is highly coupled or if the reward is not a product. 289

290 291

279

281

3.3.2 CONDITIONED JOINT FLOW NETWORK

292 As discussed training of MA-GFlowNets via training of the virtual joint GFlowNets is an approx-293 imation of the centralized training. In fact, the space of joint GFlowNets is smaller than that of the general MA-GFlowNets, as only rewards that splits into the product $R(s) = \prod_{i \in I} R^{(i)}(o^{(i)})$ 294 may be exactly sampled. If the rewards are not of this form, the training may still be subject 295 to a spurious reward or mode collapse. For instance, consider the case of $S = \{1,2\}^2$ with two 296 agents of respective positions $s_1, s_2 \in \{1, 2\}$, actions $\{(0, +1), (+1, 0), (0, 0), (+1, +1)\}$, and re-297 ward $R(s_1, s_2) = \mathbf{1}_{s_1 = s_2}$. In this case, the reward does not split and it is easy to see that inde-298 pendent agents cannot sample states proportionally to R. One may easily build more sophisticated 299 counter-examples based on this one. 300

Our proposed solution is to build a conditioned JFN inspired by augmented flows Dupont et al. 301 (2019); Huang et al. (2020) methods, which allow the bypass of architectural constraints for Nor-302 malization flows Papamakarios et al. (2021). The trick is to add a shared "hidden" state to the joint 303 MA-GFlowNets allowing the agent to synchronize. This hidden state is constant across a given 304 episode and may be understood as a cooperative strategy chosen beforehand by the agents. The 305 size of this hidden parameterization is a tradeoff: it should be large enough to allow the proper 306 parameterization of the target reward and transition but the larger the size the harder the training. 307 Formally, this simply consist in augmenting the state space and the observation spaces by a strat-308 egy space Ω to get $\tilde{S} = S \times \Omega$ and $\tilde{O}^{(i)} = O^{(i)} \times \Omega$, F_{init} is augmented by a distribution \mathbb{P} on Ω , 309 the observation projections as well as transition kernel act trivially on Ω is $T(s;\omega) = T(s)$ and 310 $p^{(i)}(s;\omega) = (p^{(i)}(s),\omega)$. The joint MA-GFlowNets theorem applies the same way, beware that 311 the observation part of $T^{(i)}$ now have a dependency on Ω even though T does not. In theory, Ω 312 may be big enough to parameterize the whole trajectory space \mathcal{T} , in which case it is possible to 313 have decoupled conditioned local transition kernels $T^{(i)}(\cdot;\omega)$ so that $\tilde{T} = T$ on a relevant domain. 314 Furthermore, the limitation on the reward is also lifted if the flow-matching property is enforced on 315 the expected joint flow $\mathbb{E}_{\omega}\mathbb{F}^{\text{joint}}$. Two possible losses may be considered: $\mathbb{E}_{\omega}\mathcal{L}_{\text{FM}}^{\text{stable}}(\mathbb{F}^{\theta,\text{joint}}(\cdot;\omega))$ or $\mathcal{L}_{\text{FM}}^{\text{stable}}(\mathbb{E}_{\omega}\mathbb{F}^{\theta,\text{joint}}(\cdot;\omega))$. The former, which we use in our experiments, is simpler to implement 316 317 but does not a priori lift the constraint on the reward.

318 The training phase of Conditioned Joint Flow Network (CJFN) is shown in Algorithm 4 in the 319 appendix. We first sample trajectories with policy

320 321 322

$$o_t^{(i)} = p_i(s_t^{(i)}) \text{ and } \pi_{\omega}^{(i)}(o_t^{(i)} \to a_t^{(i)}), \ i \in I$$
 (9)

with $a_t = (a_t^{(i)} : i \in I)$ and $s_{t+1} = T(s_t, a_t)$. Then we train the sampling policy by minimizing the FM loss $\mathbb{E}_{\omega} \mathcal{L}_{\text{FM}}^{\text{stable}}(\mathbb{F}^{\theta, \text{joint}}(\cdot; \omega))$. 323

324
 325
 326
 326
 324
 Discussion: Finally, we discuss the connection between MA-GFlowNets and multi-agent RL in Appendix C and prove some related properties.

327 328

329

4 RELATED WORKS

 Generative Flow Networks: GFlowNets is an emerging generative model that could learn a policy to generate the objects with a probability proportional to a given reward function. Nowadays, GFlowNets has achieved promising performance in many fields, such as molecule generation Bengio et al. (2021); Malkin et al. (2022); Jain et al. (2022), discrete probabilistic modeling Zhang et al. (2022), structure learning Deleu et al. (2022), domain adaptation Zhu et al. (2023), graph neural networks training Li et al. (2023b;a), and large language model training Li et al. (2023c); Hu et al. (2023); Zhang et al. (2024). This network could sample the distribution of trajectories with high rewards and can be useful in tasks where the reward distribution is more diverse.

337 GFlowNets is similar to reinforcement learning (RL) Sutton & Barto (2018). However, RL aims to 338 maximize the expected reward and often only generates the single action sequence with the highest 339 reward. Conversely, the learned policies of GFlowNets can ensure that the sampled actions are 340 proportional to the reward, making them more suitable for exploration. This exploration ability 341 makes GFlowNet promising as a new paradigm for policy optimization in the RL field, but there are 342 many problems, such as solving multi-agent collaborative tasks. Previously, the meta GFlowNets 343 algorithm Ji et al. (2024) was proposed to solve the problem of GFlowNets training under distributed 344 conditions, but it requires the observation state and task objectives of each agent to be the same, 345 which is not suitable for multi-agent problems. Later, a multi-agent GFlowNets algorithm was proposed in Luo et al. (2024), but this algorithm is an approximate algorithm without theoretical 346 support and is difficult to converge when solving large-scale multi-agent problems. In contrast, we 347 established the theory of multi-agent GFlowNets in measure space, and our algorithm can support 348 large-scale multi-agent environments, such as StarCraft missions. 349

Cooperative Multi-agent Reinforcement Learning: There exist many MARL algorithms to solve
 collaborative tasks. Two extreme algorithms for thus purpose are independent learning Tan (1993)
 and centralized training. Independent training methods regard the influence of other agents as part
 of the environment, but the team reward function often has difficulty to measure the contribution of
 each agent, resulting in the agent facing a non-stationary environment Sunehag et al. (2018); Yang
 et al. (2020).

356 On the contrary, centralized training treats the multi-agent problem as a single-agent counterpart. 357 However, this method has high combinatorial complexity and is difficult to scale beyond dozens 358 of agents Yang et al. (2019). Therefore, the most popular paradigm is centralized training and decentralized execution (CTDE), including value-based Sunehag et al. (2018); Rashid et al. (2018); 359 Son et al. (2019); Wang et al. (2020) and policy-based Lowe et al. (2017); Yu et al. (2022); Kuba 360 et al. (2022) methods. The goal of value-based methods is to decompose the joint value function 361 among the agents for decentralized execution. This requires satisfying the condition that the local 362 maximum of each agent's value function should be equal to the global maximum of the joint value 363 function. The methods, VDN Sunehag et al. (2018) and QMIX Rashid et al. (2018) employ two 364 classic and efficient factorization structures, additivity and monotonicity, respectively, despite their strict factorization method. 366

In QTRAN Son et al. (2019) and QPLEX Wang et al. (2020), extra design features are introduced for decomposition, such as the factorization structure and advantage function. The policy-based methods extend the single-agent TRPO Schulman et al. (2015) and PPO Schulman et al. (2017) into the multi-agent setting, such as MAPPO Yu et al. (2022), which has shown surprising effectiveness in cooperative multi-agent games. The goal of these algorithms is to find the policy that maximizes the long-term reward. However, it is difficult for them to learn more diverse policies in order to generate more promising results.

- 373 374
- 5 EXPERIMENTS
- 375 376

We first verify the performance of CFN on a multi-agent hyper-grid domain where partition functions can be accurately computed. We then compare the performance of CFN and CJFN with stan-



Figure 3: Mode Found and L1 error performance of different algorithms on various hyper-grid environments. Top and bottom are respectively Mode Found (higher is better) and L1 Error (lower is better). Left: Hyper-Grid v1, Middle: Hyper-Grid v2, Right: Hyper-Grid v3.

dard MCMC and some RL methods to show that our proposed sampling distributions better match normalized rewards. All our code is done using the PyTorch Paszke et al. (2019) library. We re-implemented the multi-agent RL algorithms and other baselines.

5.1 Hyper-grid Environment

396

397

402

403

404 405

406

412 413

We consider a multi-agent MDP where states are the cells of a *N*-dimensional hypercubic grid of side length *H*. In this environment, all agents start from the initialization point $x = (0, 0, \dots)$, and are only allowed to increase coordinate *i* with action a_i . In addition, each agent has a stop action. When all agents choose the stop action or reach the maximum *H* of the episode length, the entire system resets for the next round of sampling. The reward function is designed as

$$R(x) = R_0 + R_1 \prod_i \mathbb{I}(0.25 < |x_i/H - 0.5|) + R_2 \prod_i \mathbb{I}(0.3 < |x_i/H - 0.5| < 0.4), \quad (10)$$

where $x = [x_1, \dots, x_I]$ includes all agent states and the reward term $0 < R_0 \ll R_1 < R_2$ leads a distribution of modes.

416 By changing R_0 and setting it closer to 0, this environment becomes harder to solve, creating an 417 unexplored region of state space due to the sparse reward setting. We conducted experiments in 418 Hyper-grid environments with different numbers of agents and different dimensions. We used dif-419 ferent version numbers to differentiate these environments, where the higher the number is, the 420 more the number of dimensions and proxies are. The specific details about the environments and 421 experiments can be found in the appendix.

422 We compare CFN and CJFN with a modified MCMC and RL methods. In the modified MCMC 423 method Xie et al. (2021), we allow iterative reduction of coordinates on the basis of joint action 424 space and cancel the setting of stop actions to form a ergodic chain. As for the RL methods, we 425 consider the maximum entropy algorithm, i.e., multi-agent SAC Haarnoja et al. (2018), and a pre-426 vious cooperative multi-agent algorithm, i.e., MAPPO, Yu et al. (2022). Note that the maximum 427 entropy method uses the Softmax policy of the value function to make decision, so as to explore the state of other reward, which is related to our proposed algorithm. To measure the performance of 428 these methods, we use the normalized L1 error as $\mathbb{E}[|p(s_f) - \pi(s_f)| \times N]$ with $p(s_f) = R(s_f)/Z$ 429 being the sample distribution computed by the true reward, where N is cardinality of the space of 430 s_f . Moreover, we can consider the mode found theme to demonstrate the superiority of the proposed 431 algorithm.

432 Figure 3 illustrates the perfor-433 mance superiority of our pro-434 posed algorithm compared to 435 other methods in the L1 er-436 ror and Mode Found. We find that on small-scale environ-437 ments shown in Figure 3-Left, 438 CFN can achieve the best perfor-439 mance because CFN can accu-440 rately estimate the flow of joint 441 actions when the joint action 442 space dimension is small. There



Figure 4: Comparison results of JFN and Conditional JFN.

443 are two main reasons for the large 11-error index. First, we normalized the standard L1-error and 444 multiplied it by a constant to avoid the inconvenience of visualization of smaller magnitude. Sec-445 ondly, when evaluating L1-error, we only sampled 20 rounds for calculation, and with the increase 446 of the number of samples, L1-error will further decrease. As the complexity of the joint action flow that needs to be estimated increases, we find that the performance of CFN degrades. However, 447 the joint-flow based methods still achieve good estimation and maintain the speed of convergence, 448 as shown in Figure 3-Middle. Note that the RL-based methods do not achieve the expected per-449 formance. Their performance curves first rise and then fall, because as training progresses, these 450 methods tend to find the highest rewarding nodes rather than finding more patterns. Figure 4 shows 451 the performance superiority of the CJFN. When the algorithm introduces conditions to coordinate 452 multiple agents, the performance is closer to the optimal. 453

5.2 STARCRAFT

454

455

467 468

469

470

471

472

473

474

475 476

477

478 479 480

456 Figure 5 shows the performance of the proposed algorithm on the StarCraft 3m map, where (a) 457 shows the win rate comparison with different algorithms, and (b) and (c) show the decision results 458 sampled using the proposed algorithm. In the experiment, the outflow flow is calculated when the 459 flow function is large, and the maximum flow is used to calculate the win rate when sampling. It 460 can be found that since the experimental environment is not a sampling environment with diversi-461 fied rewards, although the proposed algorithm is not significantly better than other algorithms, it still illustrates its potential in large-scale decision-making. In addition, the proposed algorithm can 462 sample results with more diverse rewards, such as (b) and (c), and the number of units left implies 463 the trajectory reward. More detailed results are given in the Appendix. One thing to note is that the 464 task of the benchmark is to achieve as high a win rate as possible, which is somewhat different from 465 the goal of GFLowNets, but it can be used to verify the effectiveness of the algorithm. 466



Figure 5: The performance comparison results on the 3m map of StarCraft

6 CONCLUSION

481 482

In this paper, we discussed the policy optimization problem when GFlowNets meets the multi-agent
 systems. Different from RL, the goal of MA-GFlowNets is to find diverse samples with probability
 proportional to the reward function. Since the joint flow is equivalent to the product of independent
 flow of each agent, we designed a CTDE method to avoid the flow estimation complexity prob-

486 lem in a fully centralized algorithm and the non-stationary environment in the independent learn-487 ing process, simultaneously. Experimental results on Hyper-Grid environments and StarCraft task 488 demonstrated the superiority of the proposed algorithms.

489 Limitation and Future Work: Our theory is incomplete as it does not apply to non-cooperative 490 agents and has limited support of different game/agent terminations or initialization. A local-global 491 principle beyond independent agent policies would also be particularly interesting. Our experiments 492 do not cover the whole range of the theory in particular regarding continuous tasks and CJFN loss 493 on projected GFN. An ablation study analyzing the tradeoff of small versus big condition space Ω 494 would enlighten its importance. Finally, a metrization of the space of global GFlowNet would allow 495 a more precise functional and optimization analysis of JFN/CJFN and their limitations.

References

496 497

498

521

- 499 Emmanuel Bengio, Moksh Jain, Maksym Korablyov, Doina Precup, and Yoshua Bengio. Flow 500 network based generative models for non-iterative diverse candidate generation. In NeurIPS, 2021.
- 502 Yoshua Bengio, Salem Lahlou, Tristan Deleu, Edward J Hu, Mo Tiwari, and Emmanuel Bengio. Gflownet foundations. JMLR, 24(1):10006–10060, 2023. 504
- 505 Leo Brunswic, Yinchuan Li, Yushun Xu, Yijun Feng, Shangling Jui, and Lizhuang Ma. A theory of 506 non-acyclic generative flow networks. In AAAI Conference on Artificial Intelligence, 2024.
- 507 Lucian Busoniu, Robert Babuska, and Bart De Schutter. A comprehensive survey of multiagent rein-508 forcement learning. IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications 509 and Reviews), 38(2):156-172, 2008. 510
- 511 Lorenzo Canese, Gian Carlo Cardarilli, Luca Di Nunzio, Rocco Fazzolari, Daniele Giardino, Marco 512 Re, and Sergio Spanò. Multi-agent reinforcement learning: A review of challenges and applica-513 tions. Applied Sciences, 11(11):4948, 2021.
- 514 Tristan Deleu and Yoshua Bengio. Generative flow networks: a markov chain perspective. arXiv 515 preprint arXiv:2307.01422, 2023. 516
- 517 Tristan Deleu, António Góis, Chris Emezue, Mansi Rankawat, Simon Lacoste-Julien, Stefan Bauer, 518 and Yoshua Bengio. Bayesian structure learning with generative flow networks. In Uncertainty 519 in Artificial Intelligence, 2022.
- 520 Randal Douc, Eric Moulines, Pierre Priouret, Philippe Soulier, Randal Douc, Eric Moulines, Pierre Priouret, and Philippe Soulier. Markov chains: Basic definitions. Springer, 2018. 522
- 523 Emilien Dupont, Arnaud Doucet, and Yee Whye Teh. Augmented neural odes. Advances in neural 524 information processing systems, 32, 2019.
- Amal Feriani and Ekram Hossain. Single and multi-agent deep reinforcement learning for ai-enabled 526 wireless networks: A tutorial. IEEE Communications Surveys & Tutorials, 23(2):1226-1252, 527 2021. 528
- 529 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy 530 maximum entropy deep reinforcement learning with a stochastic actor. In ICML, 2018. 531
- Edward J Hu, Moksh Jain, Eric Elmoznino, Younesse Kaddar, Guillaume Lajoie, Yoshua Bengio, 532 and Nikolay Malkin. Amortizing intractable inference in large language models. arXiv preprint 533 arXiv:2310.04363, 2023. 534
- 535 Chin-Wei Huang, Laurent Dinh, and Aaron Courville. Augmented normalizing flows: Bridging the 536 gap between generative flows and latent variable models. arXiv preprint arXiv:2002.07101, 2020. 537
- Moksh Jain, Emmanuel Bengio, Alex Hernandez-Garcia, Jarrid Rector-Brooks, Bonaventure FP 538 Dossou, Chanakya Ajit Ekbote, Jie Fu, Tianyu Zhang, Michael Kilgour, Dinghuai Zhang, et al. 539 Biological sequence design with gflownets. In ICML, 2022.

540 Xinyuan Ji, Xu Zhang, Wei Xi, Haozhi Wang, Olga Gadyatskaya, and Yinchuan Li. Meta genera-541 tive flow networks with personalization for task-specific adaptation. *Information Sciences*, 672: 542 120569, 2024. 543 Olav Kallenberg et al. Random measures, theory and applications, volume 1. Springer, 2017. 544 Jakub Grudzien Kuba, Ruiging Chen, Munning Wen, Ying Wen, Fanglei Sun, Jun Wang, and 546 Yaodong Yang. Trust region policy optimisation in multi-agent reinforcement learning. In ICLR, 547 2022. 548 Salem Lahlou, Tristan Deleu, Pablo Lemos, Dinghuai Zhang, Alexandra Volokhova, Alex 549 Hernández-Garcia, Léna Néhale Ezzine, Yoshua Bengio, and Nikolay Malkin. A theory of con-550 tinuous generative flow networks. In ICML, 2023. 551 552 Wenqian Li, Yinchuan Li, Zhigang Li, Jianye Hao, and Yan Pang. Dag matters! gflownets enhanced 553 explainer for graph neural networks. arXiv preprint arXiv:2303.02448, 2023a. 554 Yinchuan Li, Zhigang Li, Wenqian Li, Yunfeng Shao, Yan Zheng, and Jianye Hao. Generative flow 555 networks for precise reward-oriented active learning on graphs. arXiv preprint arXiv:2304.11989, 556 2023b. 558 Yinchuan Li, Shuang Luo, Yunfeng Shao, and Jianye Hao. Gflownets with human feedback. arXiv 559 preprint arXiv:2305.07036, 2023c. 560 561 Yinchuan Li, Shuang Luo, Haozhi Wang, and Jianye Hao. Cflownets: Continuous control with generative flow networks. In ICLR, 2023d. 562 563 Ryan Lowe, Yi I Wu, Aviv Tamar, Jean Harb, OpenAI Pieter Abbeel, and Igor Mordatch. Multi-564 agent actor-critic for mixed cooperative-competitive environments. In NeurIPS, 2017. 565 566 Shuang Luo, Yinchuan Li, Shunyu Liu, Xu Zhang, Yunfeng Shao, and Chao Wu. Multi-agent 567 continuous control with generative flow networks. Neural Networks, 174:106243, 2024. 568 Nikolay Malkin, Moksh Jain, Emmanuel Bengio, Chen Sun, and Yoshua Bengio. Trajectory balance: 569 Improved credit assignment in gflownets. In NeurIPS, 2022. 570 571 Frans A Oliehoek and Christopher Amato. A concise introduction to decentralized POMDPs. 572 Springer, 2016. 573 Frans A Oliehoek, Matthijs TJ Spaan, and Nikos Vlassis. Optimal and approximate q-value func-574 tions for decentralized pomdps. Journal of Artificial Intelligence Research, 32:289–353, 2008. 575 576 George Papamakarios, Eric Nalisnick, Danilo Jimenez Rezende, Shakir Mohamed, and Balaji Lak-577 shminarayanan. Normalizing flows for probabilistic modeling and inference. Journal of Machine 578 Learning Research, 22(57):1-64, 2021. 579 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor 580 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-581 performance deep learning library. In NeurIPS, 2019. 582 583 Tabish Rashid, Mikayel Samvelyan, Christian Schroeder, Gregory Farquhar, Jakob Foerster, and 584 Shimon Whiteson. Qmix: Monotonic value function factorisation for deep multi-agent reinforce-585 ment learning. In ICML, 2018. 586 Martin Riedmiller, Roland Hafner, Thomas Lampe, Michael Neunert, Jonas Degrave, Tom Wiele, 587 Vlad Mnih, Nicolas Heess, and Jost Tobias Springenberg. Learning by playing solving sparse 588 reward tasks from scratch. In ICML, 2018. 589 John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region 591 policy optimization. In ICML, 2015. 592 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy

optimization algorithms. arXiv preprint arXiv:1707.06347, 2017.

Kyunghwan Son, Daewoo Kim, Wan Ju Kang, David Earl Hostallero, and Yung Yi. Qtran: Learning to factorize with transformation for cooperative multi-agent reinforcement learning. In <i>ICML</i> , 2019.
Matthijs TJ Spaan. Partially observable markov decision processes. In <i>Reinforcement Learning</i> , pp. 387–414. Springer, 2012.
Peter Sunehag, Guy Lever, Audrunas Gruslys, Wojciech Marian Czarnecki, Vinicius Zambaldi, Max Jaderberg, Marc Lanctot, Nicolas Sonnerat, Joel Z Leibo, Karl Tuyls, et al. Value-decomposition networks for cooperative multi-agent learning. In <i>AAMAS</i> , 2018.
Richard S Sutton and Andrew G Barto. Reinforcement learning: An introduction. MIT press, 2018.
Ming Tan. Multi-agent reinforcement learning: Independent vs. cooperative agents. In ICML, 1993.
Alexander Trott, Stephan Zheng, Caiming Xiong, and Richard Socher. Keeping your distance: Solving sparse reward tasks using self-balancing shaped rewards. In <i>NeurIPS</i> , 2019.
Jianhao Wang, Zhizhou Ren, Terry Liu, Yang Yu, and Chongjie Zhang. Qplex: Duplex dueling multi-agent q-learning. <i>arXiv preprint arXiv:2008.01062</i> , 2020.
Yutong Xie, Chence Shi, Hao Zhou, Yuwei Yang, Weinan Zhang, Yong Yu, and Lei Li. {MARS}: Markov molecular sampling for multi-objective drug discovery. In <i>ICLR</i> , 2021.
Yaodong Yang, Rasul Tutunov, Phu Sakulwongtana, Haitham Bou Ammar, and Jun Wang. α^{α} -rank: Scalable multi-agent evaluation through evolution. 2019.
Yaodong Yang, Ying Wen, Jun Wang, Liheng Chen, Kun Shao, David Mguni, and Weinan Zhang. Multi-agent determinantal q-learning. In <i>ICML</i> , 2020.
Chao Yu, Akash Velu, Eugene Vinitsky, Yu Wang, Alexandre Bayen, and Yi Wu. The surprising effectiveness of ppo in cooperative, multi-agent games. In <i>NeurIPS</i> , 2022.
Dinghuai Zhang, Nikolay Malkin, Zhen Liu, Alexandra Volokhova, Aaron Courville, and Yoshua Bengio. Generative flow networks for discrete probabilistic modeling. In <i>ICML</i> , 2022.
Dinghuai Zhang, Yizhe Zhang, Jiatao Gu, Ruixiang Zhang, Josh Susskind, Navdeep Jaitly, and Shuangfei Zhai. Improving gflownets for text-to-image diffusion alignment. <i>arXiv preprint arXiv:2406.00633</i> , 2024.
Kaiqing Zhang, Zhuoran Yang, and Tamer Başar. Multi-agent reinforcement learning: A selective overview of theories and algorithms. <i>Handbook of Reinforcement Learning and Control</i> , pp. 321–384, 2021.
Didi Zhu, Yinchuan Li, Yunfeng Shao, Jianye Hao, Fei Wu, Kun Kuang, Jun Xiao, and Chao Wu. Generalized universal domain adaptation with generative flow networks. In <i>Proceedings of the</i> 31st ACM International Conference on Multimedia, pp. 8304–8315, 2023.

648 JOINT FLOW THEORY А 649

The goal of this section is to lay down so elementary points on the measurable theory of MA-GFlowNets as well as prove the main theorem on the joint GFlowNet.

A.1 NOTATIONS ON MEASURES AND KERNELS

655 We mostly use notations from Douc et al. (2018) regarding kernels and measures. In the whole section, since we deal with technicalities, we use kernel type notations for image by kernels and 656 maps (seen as deterministic kernels). So that for a kernel $K : X \to Y$ and a measure μ on X 657 we denote by μK the measure on Y defined by $\mu K(B) = \int_{x \in X} K(x \to B) d\mu(x)$ for $B \subset Y$ 658 measurable and $\mu \otimes K$ is the measure on $X \times Y$ so that $\mu \otimes K(A \times B) = \int_{x \in A} K(x \to B) d\mu(x)$. 659 Recall that a measure ν dominates a measure μ which is denoted $\mu \ll \nu$, if for all measurable A, 660 $\nu(A) = 0 \Rightarrow \mu(A) = 0$. The Radon-Nykodim Theorem ensures that if $\mu \ll \nu$ and μ, ν are finite then 661 there exists $\varphi \in L^1(\nu)$ so that $\mu = \varphi \nu$. This function φ is called the Radon-Nykodim derivative and 662 is denoted $\frac{d\mu}{d\nu}$. We favor notations $\mu(A \to B)$ when μ is a measure on $X \times Y$ and $A \subset X$ and $B \subset Y$; 663 also $\mu(A \rightarrow \cdot)$ means the measure $B \mapsto \mu(A \rightarrow B)$. 664

665 666

667

668

669

672 673

674 675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

697

650

651

652 653

654

A.2 AN INTRODUCTION FOR NOTATIONS

We understand that our formalism is abstract, this section is devoted justifying our choices and providing examples.

670 A.2.1 MOTIVATIONS 671

> To begin with, our motivation to formalize the action space as a measurable bundle $\mathcal{A} := \{(s, a) \mid s \in \{a, b\} | s \in \{a, b\} \}$ $\mathcal{S}, a \in \mathcal{A}_s \} \xrightarrow{S} \mathcal{S}$ is three fold:

- 1. The available actions from a state may depend on the state itself: on a grid, the actions available while on the boundary of the grid are certainly more limited than while in the middle. More generally, on a graph, actions are typically formalized by edges $s \xrightarrow{a} s'$ of the graph, the data of an edge contains both the origin s and the destination s'. In other words, on graphs, actions are bundled with an origin state. It is thus natural to consider the actions as bundled with the origin state. When an agent is transiting from a state to another via an action, the state map tells where it comes from while the transition map tells where it is going.
- 2. We want our formalism to cover as many cases as possible in a unified way: Graphs, vector spaces with linear group actions or mixture of discrete and continuous state spaces. Measures and measurable spaces provide a nice formalism to capture the quantity of reward on a given set or a probability distribution.
 - 3. We want a well-founded and possibly standardized mathematical formalism. In particular, the policy takes as input a state and returns a distribution of actions. the actions should correspond to the input state. To avoid having a cumbersome notion of policy as a family of distributions π_s each on \mathcal{A}_s , we prefer to consider the union of the state-dependent action spaces $\mathcal{A} := \bigcup_{s \in S} \mathcal{A}_s$ and define the policy as Markov kernel $S \to \mathcal{A}$. However, we still need to satisfy the constraint that the distribution $\pi(s)$ is supported by \mathcal{A}_s . Bundles are usual mathemcatical objects formalizing such situations and constraints and are thus well suited for this purpose and the constraint is easily expressed as $S \circ \pi(s) = s, \forall s \in S$.
- 694 696
 - Our synthetic formalism comes with a few drawbacks due to the level of abstraction:
 - 1. The notation $\pi(s)$ differs from the more common notation $\pi(s, a)$ as the action already contains s implicitly.
- 2. We need to use Radon-Nikodym derivative. At a given state, on a graph, a GFlowNets has 699 a probability to stop 700 $\mathbb{P}(\mathrm{STOP}|s) = \frac{R(s)}{F_{\mathrm{out}}(s)}.$

On a continuous statespace with reference measure λ , the stop probability is

$$\mathbb{P}(\mathrm{STOP}|s) = \frac{r(s)}{f_{\mathrm{out}}(s)}$$

where r(s) is the density of reward at s and $f_{out}(s)$ is the density of outflow at s. A natural measure-theoretic way of writing these equations as one is via Radon-Nikodym derivation: given two measures μ, ν ; if $\mu(X) = 0 \Rightarrow \nu(X) = 0$ for any measurable $X \subset S$ then μ is said to dominate ν and, by Radon-Nikodym Theorem, there exists a measurable function $\varphi \in L^1(\mu)$ such that $\nu(X) = \int_{x \in X} \varphi(x) d\nu(x)$ for all measurable $X \subset S$. This φ is the Radon-Nikodym derivative $\frac{d\nu}{d\mu}$.

If one has a measure λ dominating both R and $F_{\rm out}$ and if $F_{\rm out}$ dominated R then

$$\mathbb{P}(\mathrm{STOP}|s) \coloneqq \frac{dR}{dF_{\mathrm{out}}}(s) = \frac{dR}{d\lambda}(s) \times \left(\frac{dF_{\mathrm{out}}}{d\lambda}\right)^{-1}$$

When S is discrete, we choose λ as the counting measure and we recover the graph formula above. When S is continuous, we choose λ as the Lebesgue measure and we recover the second formula.

A.2.2 EXAMPLE 720

702

704

706

708

709

710 711

716

717

718 719

729

739 740

721 Consider the D-dimensional W-width hypergrid case with agent set I, see Figure 6. The state space 722 is the finite set $S = (\{1, \dots, W\}^D)^I$, say each agent only observes its own position on the grid so 723 that $\mathcal{O}^{(i)} = \{1, \dots, W\}^D$. the observation-dependent action space of the *i*-th agent $\mathcal{A}_{o^{(i)}}^{(i)}$ is a subset 724 of $H := \{\pm \mathbf{1}_k : 1 \le k \le W\}$ where $\mathbf{1}_k$ is the hot-one array $(0, \dots, 0, 1, 0, \dots, 0)$ with a one at the k-th 725 coordinate. The set $\mathcal{A}_{\alpha^{(i)}}^{(i)}$ depends on s: if $1 < s_k < W$ then $\mathcal{A}_{\alpha^{(i)}}^{(i)} = \{\pm \mathbf{1}_k : 1 \le k \le W\} \cup \{\text{STOP}\}$ 726 727 but if $s_k = 1$ then $-\mathbf{1}_k \notin \mathcal{A}_{\alpha(i)}^{(i)}$ and similarly if $s_k = W$. The local total action space is then 728

$$\mathcal{A}_{o^{(i)}}^{(i)} = \{(s,a) \mid 1 \le s_k \le W \text{ and } 1 \le s_k + a_k \le W\} \cup \{\text{STOP}\} \subset \{1, \dots, W\}^D \times H \cup \{\text{STOP}\}.$$

730 The local state maps $S^{(i)}$ is $S^{(i)}(o^{(i)}, a) = o^{(i)}$. Since each agent may choose its action freely, for 731 any $s \in S$, $A_s = \prod_{i \in I} A_{o^{(i)}} / \sim$ however, since $A_{o^{(i)}}$ depends on i and s then $A \neq \prod_{i \in I} A_{o^{(i)}} / \sim$. 732

The local transition kernel $T^{(i)}$ depends both on the global transition kernel and the policies of all 733 734 the agents. Two possible choices of transitions depend on whether the agent interacts or not. In the non-interacting case $T_1(s, a) = s + a$. If agents may not occupy the same position then the transition 735 rejects the action if the agent moving would put them in the same position; so $T_2(s, a) = s + a$ if 736 s + a is legal, otherwise $p^{(i)} \circ T_2(s, a) = o^{(i)}$ for some *i*. The simplest T_2 is to choose $T_2(s, a) = s$ 737 if s + a is illegal. In this case 738

$$T_2^{(i)}(o^{(i)}, a^{(i)}) = \mathbb{P}(s + a \text{ is legal}|o^{(i)}, a^{(i)})\delta_{o^{(i)} + a^{(i)}} + \mathbb{P}(s + a \text{ is illegal}|o^{(i)}, a^{(i)})\delta_{o^{(i)}}$$

Clearly, $\mathbb{P}(s + a \text{ is legal}|o^{(i)}, a^{(i)})$ depends on the policies and positions of all the agents, then so 741 does the local transition kernels $T_2^{(i)}$. 742

A non-negative measure μ on S is any function of the form $\mu(X) = \sum_{x \in X} f(x)$ with $f : S \to \mathbb{R}_+$ any function. Defining the counting measure $\lambda(X) \coloneqq \sum_{x \in X} 1 = \operatorname{Card}(X)$ we have $\mu = f\lambda$ as measures on S, or equivalently, $\frac{d\mu}{d\lambda} = f$. We may thus translate any reward or probability distribution on such a hypergrid as a measure. 743 744 745 746 747

A policy is a Markov kernel $\mathcal{S} \to \mathcal{A}$ such that $\mathcal{S} \circ \pi = \text{Identity}$. More concretely, it means we have 748 a function that associates to any state s a probability distribution on \mathcal{A} with support on elements of 749 the form (s, a) with $a \in A_s$. From the description of measures, such a policy is fully described by a 750 function $q : \mathcal{A} \to \mathbb{R}_+$ such that 751 $\forall s \in \mathcal{S}, \sum_{a \in \mathcal{A}_s} q(s, a) = 1.$

753 The policy is then $\pi(s) = \sum_{a \in \mathcal{A}_s} q(s, a) \delta_{(s,a)}$. 754

A GFlowNet on this hypergrid in reward-less notations is given by $(F_{\text{init}}, \pi^*, F_{\text{out}}^*)$. Now, F_{init} is 755 any measure on S, it may be given by a pre-chosen family of categorical distribution of the finite set


$$\in \mathcal{S}, \quad \mathcal{A}_s \smallsetminus \{\text{STOP}\} = \prod_{i \in I} \left(\mathcal{A}_{p^{(i)}(s)}^{(i)} \smallsetminus \{\text{STOP}\} \right).$$

 $\forall s$

• A multi-agent interactive environment is a tuple $(S, A, S, T, \mathcal{O}^{(i)}, \mathcal{A}^{(i)}, S^{(i)}, p^{(i)})_{i \in I}$ where $(S, A, S, \mathcal{O}^{(i)}, \mathcal{A}^{(i)}, S^{(i)}, p^{(i)})_{i \in I}$ is a multi-agent action environment and (S, A, S, T) is a mono-agent interactive environment.

• A multi-agent game environment is a tuple $(S, A, S, T, R, \mathcal{O}^{(i)}, \mathcal{A}^{(i)}, S^{(i)}, p^{(i)})_{i \in I}$ such that $(S, A, S, T, \mathcal{O}^{(i)}, \mathcal{A}^{(i)}, S^{(i)}, p^{(i)})_{i \in I}$ is multi-agent interactive environment and (S, A, S, T, R) is a mono-agent game environment.

818 A.4 GFLOWNET IN A GAME ENVIRONMENT

A generative flow networks may be formally defined on an action environment (S, A, S), as a triple $(\pi^*, F_{out}^*, F_{init})$ where $\pi^* : S \to A$ is a Markov kernel such that $\pi^*S = Id_S$, F_{out}^* and F_{init} are a finite non-negative measures on S. Furthermore, we assume that for all $s \in S, \pi^*(s \to STOP_s) = 0$.

On an interactive environment (S, A, S, T), given a GFlowNet $(\pi^*, F_{out}^*, F_{init})$, we define the ongoing flow as $F_{in} \coloneqq F_{out}^* \pi^* T + F_{init}$ and the GFlowNet induces an virtual reward $\hat{R} \coloneqq F_{in} - F_{out}^*$. Note that the virtual reward is always finite as the star-outflow and the initial flow are both finite and π^* and T are Markovian.

Definition 1 (Weak Flow-Matching Constraint) The weak flow-matching constraint is defined as

 $\hat{R} \ge 0 \tag{11}$

If the GFlowNet satisfies the weak flow-matching constraint, we may define a virtual GFlowNet policy as

$$\hat{\pi} \coloneqq \frac{dF_{\text{out}}^*}{dF_{\text{in}}} \pi^* \tag{12}$$

where δ_{STOP} is the deterministic Markov kernel sending any s to STOP_s . The virtual action and edge flows are:

$$\hat{F}_{action} \coloneqq F_{in} \otimes \hat{\pi} \in \mathcal{M}^+(\mathcal{S} \times \mathcal{A}); \tag{13}$$

$$\hat{F}_{edge} \coloneqq F_{in} \otimes (\hat{\pi}T) \in \mathcal{M}^+(\mathcal{S} \times \mathcal{S}).$$
(14)

In a game environment, a GFlowNet comes with an outgoing flow, a natural policy, a natural action
 flow and a natural edge flow

$$F_{\text{out}} \coloneqq F_{\text{out}}^* + R \tag{15}$$

$$\pi \coloneqq \frac{dF_{\text{out}}^*}{dF_{\text{out}}} \pi^* \tag{16}$$

$$F_{\text{edge}} \coloneqq F_{\text{out}} \otimes (\pi T) \in \mathcal{M}^+(\mathcal{S} \times \mathcal{S})$$
(17)

$$F_{\text{action}} \coloneqq F_{\text{out}} \otimes \pi \in \mathcal{M}^+(\mathcal{S} \times \mathcal{A}).$$
(18)

By abuse of notation we also write F_{action} (resp. F_{action}) for $F_{\text{out}}\pi$ (resp. $F_{\text{in}}\pi$). and the flow-matching property may be rewritten as follows.

Definition 2 (Flow-Matching Constraint) The flow-matching constraint on a Game environment (S, A, S, T, R) is defined as

$$\hat{R} = \mathbb{E}(R). \tag{19}$$

Remark 1 In an interactive environment $(S, A, S, T, O^{(i)}, A^{(i)}, S^{(i)}, p^{(i)})_{i \in I}$, a GFlowNet satisfy-856 ing the weak flow-matching constraint satisfies the (strong) flow-matching constraint on the Game 857 environment $(S, A, S, T, \hat{R}, O^{(i)}, A^{(i)}, S^{(i)}, p^{(i)})_{i \in I}$.

We may recover part of the GFlowNets $(\pi^*, F_{out}^*, F_{init})$ from any of F_{action} , \hat{F}_{action} as in general:

$$\pi^*(x \to A) = \frac{dF_{action}(\cdot \to A \smallsetminus \text{STOP})}{dF_{action}(\cdot \to A \smallsetminus \text{STOP})} = \frac{d\hat{F}_{action}(\cdot \to A \smallsetminus \text{STOP})}{d\hat{F}_{action}(\cdot \to A \smallsetminus \text{STOP})}$$
(20)

$$R = F_{action}(\cdot \to \text{STOP})$$
 $\hat{R} = \hat{F}_{action}(\cdot \to \text{STOP})$ (21)

$$F_{\text{out}}^* = F_{\text{action}}(\cdot \to \mathcal{A}) - R = \hat{F}_{\text{action}}(\cdot \to \mathcal{A}) - \hat{R}$$
(22)

$$F_{\text{init}} = F_{\text{out}}^* T + \hat{R} \tag{23}$$

867 If the flow-matching constraint is satisfied, then

$$F_{\text{init}} = F_{\text{out}}^* T + R. \tag{24}$$

Before going further, the presence densities.

Definition 3 Let $\mathbb{F} = (\pi^*, F_{out}, F_{init})$ be a GFlowNet in an interactive environment $(\mathcal{S}, \mathcal{A}, S, T, \mathcal{O}^{(i)}, \mathcal{A}^{(i)}, S^{(i)}, p^{(i)})_{i \in I}$.

The initial density of \mathbb{F} is the probability distribution

$$\nu_{\mathbb{F},\text{init}} \coloneqq \frac{1}{F_{\text{init}}(\mathcal{S})} F_{\text{init}}$$

The virtual presence density of \mathbb{F} is the probability distribution $\hat{\nu}_{\mathbb{F}}$ defined by

$$\hat{\nu}_{\mathbb{F}} \propto \sum_{t=0}^{\infty} \nu_{\mathbb{F},\text{init}} \hat{\pi}^t.$$

The anticipated presence density of \mathbb{F} is the probability distribution $\overline{\nu}_{\mathbb{F}}$ defined by

$$\overline{\nu}_{\mathbb{F}} \coloneqq \frac{1}{F_{\mathrm{in}}(\mathcal{S})} F_{\mathrm{in}}.$$

In a game environment, the presence density of \mathbb{F} is the probability distribution $\nu_{\mathbb{F}}$ defined by

$$u_{\mathbb{F}} \propto \sum_{t=0}^{\infty} \nu_{\mathbb{F}, \text{init}} \pi^t.$$

Lemma 1 Let \mathbb{F} be a GFlowNet in an interactive environment satisfying the weak flow-matching constraint. If $\hat{\nu}_{\mathbb{F}} \gg \overline{\nu}_{\mathbb{F}}$, then $\hat{\nu}_{\mathbb{F}} = \overline{\nu}_{\mathbb{F}}$.

Proof 1 Let $(S, A, S, T, \mathcal{O}^{(i)}, \mathcal{A}^{(i)}, S^{(i)}, p^{(i)})_{i \in I}$ be the interactive environment and let \mathbb{F} = $(\pi^*, F_{out}, F_{init})$. To begin with, $\mathbb{F}' := (\pi^*, F_{init}(S)\hat{\nu}_{\mathbb{F}} - \hat{R}, F_{init})$ is a GFlowNet satisfying the strong flow-matching constraint for reward \hat{R} , its edgeflow F'_{edge} may be compared to the edgeflow F_{edge} of \mathbb{F} : by Proposition 2 of Brunswic et al. (2024), we have $F_{edge} \ge F'_{edge}$, and the difference $F_{edge} - F'_{edge}$ is a 0-flow in the sense this same article. Also, the domination hypothesis implies that $F'_{edge} \gg F^0_{edge} \coloneqq F_{edge} - F'_{edge}$. Since the edge-policy of F_{edge} is the same as that of F'_{edge} . we deduce that it is also the same as F_{edge}^0 . By the same Proposition 2, we have $F'_{out}\pi^t \xrightarrow{t \to +\infty} 0$, therefore, $\mu \pi^t \xrightarrow{t \to +\infty} 0$ for any $\mu \ll F'_{out}$. Again by domination, $F'_{edge} \gg F^0_{edge}$ we deduce that $F'_{\text{out}} \gg F^0_{\text{out}}$. Therefore, $F^0_{\text{out}}\pi^t \xrightarrow{t \to +\infty} 0$. Finally, since 0 is a 0-flow, $F^0_{\text{out}}\pi = F^0_{\text{out}}$, we deduce that $F^0_{\text{out}} = 0$ and thus $F_{\text{edge}} = F'_{\text{edge}}$ ie $\hat{\nu}_{\mathbb{F}} = \overline{\nu}_{\mathbb{F}}$.

Remark 2 As long as the GFlowNets considered are trained using an FM-loss on a training training distribution ν_{state} extracted from trajectory distributions $\hat{\nu}_{\mathbb{F}}$ or $\nu_{\mathbb{F}}$ of the GFlowNets themselves, we may assume that $\hat{\nu}_{\mathbb{F}} \gg \overline{\nu}_{\mathbb{F}}$ as flows are only evaluated on a distribution dominated by $\nu_{\mathbb{F}}$. The singular part with respect to $\nu_{\mathbb{F}}$ is irrelevant for training purposes as well as inference purposes. Therefore, we may generally assume that $\hat{\nu} = \overline{\nu}$

913
914 Remark 3 The main interest of the virtual reward R̂ is for cases where the reward is not accessible
915 or expensive to compute. Since a GFlowNet satisfying the weak flow-matching property always
916 satisfies the strong flow-matching property for the reward R̂, the sampling Theorem usually applies
917 to R̂. Therefore, R̂ may be used as a reward during inference instead of the true reward R so that we actually sample using the policy π̂ instead of π.

918 A.5 MA-GFLOWNETS IN MULTI-AGENT ENVIRONMENTS (I): PRELIMINARIES

⁹²⁰ To begin with, let us define a MA-GFlowNet on a multi-agent environment.

Definition 4 An MA-GFlowNet on a multi-agent action environment is the data of a global GFlowNet \mathbb{F} on $(\mathcal{S}, \mathcal{A}, S)$ and a collection of local GFlowNets $\mathbb{F}^{(i)}$ on $(\mathcal{O}^{(i)}, \mathcal{A}^{(i)}, S^{(i)})$ for $i \in I$.

We give ourselves a multi-agent interactive environment $(S, A, S, T, O^{(i)}, A^{(i)}, S^{(i)}, p^{(i)})$. We wish to clarify the links between local and global GFlowNet.

- A priori, there the local GFlowNets are merely defined on an action environment, they lack both the local transition kernel $T^{(i)}$ and the reward $R^{(i)}$.
- Given a global GFlowNet, we wish to define local GFlowNets.
- Given a family of local GFlowNets, we wish to define a global GFlowNet.

For simplicity sake, for any GFlowNet \mathbb{F} defined on an interactive environment satisfying the weak flow-matching constraint, we set $R = \hat{R}$ and apply remark 2 assume that $\hat{\nu}_{\mathbb{F}} = \overline{\nu}_{\mathbb{F}} = \nu_{\mathbb{F}}$.

Definition 5 Let $(S, A, S, T, O^{(i)}, A^{(i)}, S^{(i)}, p^{(i)})$ be a multi-agent interactive environment and let $\mathbb{F} = (\pi^*, F_{out}^*, F_{init})$ be a GFlowNet on (S, A) satisfying the weak flow-matching constraint. We introduce the following:

- the local presence probability distribution $\nu_{\mathbb{R}}^{(i)} \coloneqq \nu_{\mathbb{R}} p^{(i)}$;
- the measure map $o^{(i)} \mapsto \nu_{\mathbb{F}|o^{(i)}}$ as the disintegration of $\nu_{\mathbb{F}}$ by $p^{(i)}$
- the Markov kernel $\tilde{\pi}^{(i)} : \mathcal{O}^{(i)} \to \mathcal{A}$ by $\delta_{o^{(i)}} \tilde{\pi}^{(i)} := \nu_{\mathbb{F}|o^{(i)}} \pi$;
- the Markov kernel $\pi^{(i)}: \mathcal{O}^{(i)} \to \mathcal{A}^{(i)}$ by $\pi^{(i)} = \tilde{\pi}^{(i)} p^{(i)}$;
- the Markov kernel $T^{(i)}: \mathcal{A}^{(i)} \to \mathcal{O}^{(i)}$ by $T^{(i)} = S^{(i)} \tilde{\pi}^{(i)} T p^{(i)}$;

The situation may be summarized by the following diagram:



Before going further, we need to check that these definitions are somewhat consistent.

Lemma 2 The following diagrams are commutative in the category of probability spaces.



Proof 2 For the left diagram, with the definition chosen, we only need to check that $\nu_{\mathbb{F}}^{(i)} \tilde{\pi}^{(i)} = \nu_{\mathbb{F}} \pi$. For all $\varphi \in L^1(\mathcal{A}, \nu_{\mathbb{F}} \pi)$ we have

$$\begin{split} \int_{s \in \mathcal{A}} \varphi(a) d(\nu_{\mathbb{F}} \pi)(a) &= \int_{s \in \mathcal{S}} \int_{a \in \mathcal{A}} \varphi(a) d\pi(s, a) d\nu_{\mathbb{F}}(s) \\ &= \int_{o^{(i)} \in \mathcal{O}^{(i)}} \int_{s \in (p^{(i)})^{-1}(o^{(i)})} \int_{a \in \mathcal{A}} \varphi(a) d\pi(s, a) d\nu_{\mathbb{F}}|_{o^{(i)}}(s) d\nu_{\mathbb{F}}^{(i)}(o^{(i)}) \\ &= \int_{o^{(i)} \in \mathcal{O}^{(i)}} \int_{a \in \mathcal{A}} \varphi(a) d\tilde{\pi}^{(i)}(a) d\nu_{\mathbb{F}}^{(i)}(o^{(i)}) \\ &= \int_{a \in \mathcal{A}} \varphi(a) d(\nu_{\mathbb{F}}^{(i)} \tilde{\pi}^{(i)})(a). \end{split}$$

> For the right diagram, we need to check that $\nu_{\mathbb{F}}\pi p^{(i)} = \nu^{(i)}\pi^{(i)}$ and that $\nu_{\mathbb{F}}\pi T p^{(i)} = \nu_{\mathbb{F}}^{(i)}\pi^{(i)}T^{(i)}$. We already proved the first equality for the left diagram and for the second:

$$\nu_{\mathbb{F}} \pi p^{(i)} T^{(i)} \coloneqq \nu_{\mathbb{F}} \underbrace{\pi p^{(i)} S^{(i)}}_{=p^{(i)}} \tilde{\pi}^{(i)} T p^{(i)} = \underbrace{\nu_{\mathbb{F}} p^{(i)}}_{\nu_{\mathbb{F}}^{(i)}} \tilde{\pi}^{(i)} T p^{(i)} = \nu_{\mathbb{F}}^{(i)} \pi^{(i)} T^{(i)}$$

We see that from a global GFlowNet, we may build local policies as well as local transition kernels. These policies and transitions are natural in the sense that of local the induced local agent policy an transition are exactly the one wed would have if the observations of the other agents were provided as a random external parameter. The local rewards are then stochastics depending on the state of the global GFlowNet.

A.6 MA-GFLOWNETS IN MULTI-AGENT ENVIRONMENTS (II): FROM LOCAL TO GLOBAL

We would like to settle construction of global GFlowNet from local ones, key difficulties arise:

- the global distributions induce local ones but the coupling of the local distributions may be non trivial;
- the defining the star-outflow and initial flow requires to find proportionality constants

$$F_{\mathrm{in}}(\mathcal{O}^{(i)}) \propto \nu_{\mathbb{F}}^{(i)} \qquad F_{\mathrm{init}}^{(i)} \propto \nu_{\mathbb{F}^{(i)},\mathrm{init}};$$

• The coupling of the local transition kernels $T^{(i)}$ and the global one is in general non-trivial.

We try to solve these issues by looking at the simplest coupling: independent local agents. Recall that $\mathcal{A}_s^* = \prod_{i \in I} \mathcal{A}_s^{(i),*}$ therefore, independent coupling means that $\pi^*(s \to \cdot) = \prod_{i \in I} \pi^{(i),*}(o^{(i)} \to \cdot)$. We may generalize this relation to a coupling of GFlowNets writing $F_{\text{action}}(\prod_{i \in I} O^{(i)} \to \prod_{i \in I} \mathcal{A}_{\text{action}}^{(i)}) = \prod_{i \in I} F_{\text{action}}^{(i)}(O^{(i)} \to \mathcal{A}^{(i)})$. We are led to following the definition:

1011 Definition 6 Let $(S, A, S, T, O^{(i)}, A^{(i)}, S^{(i)}, p^{(i)})$ be a multi-agent interactive environment and let **1012** $\mathbb{F} = (\pi^*, F^*_{out}, F_{init})$ be a global GFlowNet on it satisfying the weak flow-matching constraint. The **1013** GFlowNet \mathbb{F} is said to be

• star-split if for some local GFlowNets $\mathbb{F}^{(i)}$ and $\forall A^{(i)} \subset \mathcal{A}^{(i)} \setminus \text{STOP}$ we have:

$$F_{\text{action}}(\prod_{i \in I} A^{(i)}) = \prod_{i \in I} F_{\text{action}}^{(i)}(A^{(i)}).$$
(25)

• strongly star-split if for some local GFlowNets $\mathbb{F}^{(i)}$ and $\forall A^{(i)}, B^{(i)} \subset \mathcal{O}^{(i)}$ we have:

$$F_{\text{edge}}\left(\prod_{i \in I} A^{(i)} \to \prod_{i \in I} B^{(i)}\right) = \prod_{i \in I} F_{\text{edge}}^{(i)}(A^{(i)} \to B^{(i)}).$$

$$(26)$$

1023 The local GFlowNets $\mathbb{F}^{(i)}$ are called the components of the global GFlowNet \mathbb{F} .

1025 However we encounter an additional difficulty: what happens when an agent decides to stop the game ? Indeed, local agents have their own STOP action, we then have at least three behaviors.

Unilateral Stop: if any agent decides to stop, the game stops and reward is awarded.
 Asynchronous Unanimous Stop: if an agent decides to stop, it does not act anymore, waits

- for the other to leave the game and then reward is awarded only when all agents stopped.
- 3. Synchronous Unanimous Stop: if an agent decides to stop but some other does not, then the stop action is rejected and the agent plays a non-stopping action.

Similar variations may be considered for how the initialization of agents:

- 1. Asynchronous Start: the game has a free number of player, agents may enter the game while other are already playing.
- 2. Synchronous Start: the game has a fixed number of players, and agents all start at the same time.

These 6 possible combinaisons leads to slight variations on the formalization of MA-GFlowNetsfrom local GFlowNets.

1041

1047 1048

1051

1053 1054 1055

1059

1062 1063 1064

1076 1077 1078

1026

1029

1030

1031

1034

1035

1036

1037

1042 A.7 INITIAL LOCAL-GLOBAL CONSISTENCIES 1043

1044 Let us formalize Asynchronous and Synchronous starts. In synchronous case, the agents are ini-1045 tially distributed according to their own initial distributions and independently. Therefore, ν_{init} is a 1046 product and

$$F_{\text{init}} \propto \nu_{\text{init}} = \prod_{i \in I} \nu_{\text{init}}^{(i)} \propto \prod_{i \in I} F_{\text{init}}^{(i)}$$

1049 Also, by strong star-splitting property, $F_{in}^* = \prod_{i \in I} F_{in}^{(i),*}$. By $F_{in} = F_{init} + F_{in}^*$ we obtain the definition below.

Definition 7 A strongly star-split global GFlowNet is said to have Synchronous start if

$$F_{\text{in}} = \prod_{i \in I} F_{\text{init}}^{(i)} + \prod_{i \in I} F_{\text{in}}^{(i)}$$

On the other hand, in the asynchronous case, an incoming agent may "bind" to agent arriving at the same time and other already there hence, the initial flow is a combination of any of the products

$$F_{\text{init}} = \sum \prod_{i \in \{\text{incoming}\}} F_{\text{init}}^{(i)} \prod_{j \in \{\text{already in}\}} F_{\text{init}}^{(j),*} = \prod_{i \in I} (F_{\text{init}}^{(i)} + F_{\text{in}}^{(i),*}) - \prod_{i \in I} F_{\text{in}}^{(i),*}.$$

Definition 8 A strongly star-split global GFlowNet is said to have Asynchronous start if

$$F_{\rm in} = \prod_{i \in I} (F_{\rm init}^{(i)} + F_{\rm in}^{(i),*})$$

A.8 TERMINAL LOCAL-GLOBAL CONSISTENCIES

We focus on terminal behaviors 1 and 2 which we formalize as follows. Local-global consistency consists in describing the formal structure linking local environments with global ones. The product structure of the action space is slightly different depending on the terminal behavior. It happens that we may up to formalization, we may cast Asynchronous Unanimous STOP as a particular case of Unilateral STOP local-global consistency. More precisely:

Definition 9 (Unilateral STOP Local-Global Consistency) With the same notations as above, we
 say that a multi-agent action environment has unilateral STOP if

$$\mathcal{A}_{s} \coloneqq \left(\prod_{i \in I} \mathcal{A}_{o^{(i)}}\right) / \sim \qquad a_{1} \sim a_{2} \Leftrightarrow \exists i, j \in I, a_{1}^{(i)} = \operatorname{STOP}^{(i)}, a_{2}^{(j)} = \operatorname{STOP}^{(j)}.$$
(27)

Definition 10 (Asynchronous Unanimous STOP Local-Global Consistency) With the same notations as above, we say that a multi-agent game environment has Asynchronous Unanimous

STOP if is has Unilateral STOP and the observation space $\mathcal{O}^{(i)}$ may be decomposed into $\mathcal{O}^{(i)}$ = 1081 $\mathcal{O}_{life}^{(i)} \cup \mathcal{O}_{purgatory}^{(i)}$ and for any observation $o^{(i)} \in \mathcal{O}_{life}^{(i)}$ we have some $\tilde{o}^{(i)} \in \mathcal{O}_{purgatory}^{(i)}$ such that : 1082

1084

1088 where the value on top of arrows are constrained flow values. 1089

1090 The formal definition of Unilateral STOP is straightforward as any local STOP activates the global 1091 STOP so that any combination of local actions that contains at least one STOP is actually a global STOP. The quotient by the equivalence relation formalizes this property. Regarding Asynchronous 1093 Unanimous STOP, the chosen formalization allows to store the last observation of an agent while it 1094 is put on hold until global STOP. Indeed, a standard action (\neq STOP) is invoked to enter purgatory, 1095 the reward is supported on purgatory and as long as all the agent are not in purgatory its value is zero (recall that from the viewpoint of a given agent, $R^{(i)}$ is stochastic but in fact depends on the whole 1096 global state). The local STOP action is then never technically called on an "alive" observation, once 1097 in purgatory the ε self-transition is called by default as long as the reward is non zero, hence until all agents are in purgatory. When the reward is activated, the policy at a purgatory state $\tilde{o}^{(i)}$ is then 1099 $\frac{d\varepsilon}{d(\varepsilon+R^{(i)})}\delta_{\tilde{o}^{(i)}} + \frac{dR^{(i)}}{d(\varepsilon+R^{(i)})}\delta_{\text{STOP}}. \text{ As } \varepsilon \to 0^+, \text{ the policy becomes equivalent to "if reward then STOP,"}$ 1100 else WAIT". This behavior is exactly the informal description of Asynchronous Unanimous STOP, 1101 1102 the formalization is rather arbitrary and does not limit the applicability as it simply helps deriving 1103 formulas more easily.

1104 We now prove Theorem 2. 1105

Theorem 4 Let $(S, A, S, T, \mathcal{O}^{(i)}, \mathcal{A}^{(i)}, S^{(i)}, p^{(i)})$ be a multi-agent interactive environment. Let 1106 $\mathbb{F}^{(i)}$ be non-zero GFlowNets on $(\mathcal{O}^{(i)}, \mathcal{A}^{(i)}, S^{(i)})$ for $i \in I$ satisfying the weak flow-1107 1108 matching constraint, then there exists a transition kernel \tilde{T} and a star-split GFlowNet on 1109 $(\mathcal{S}, \mathcal{A}, S, \tilde{T}, \mathcal{O}^{(i)}, \mathcal{A}^{(i)}, S^{(i)}, p^{(i)})$ whose components are the $\mathbb{F}^{(i)}$.

Furthermore, 1111

1110

1112

1113

1114

1115 1116

1117

1118

1119

1128

• if the multi-agent environment is a game environment with Asynchronous Unanimous STOP and if the global GFlowNet satisfies the strong flow-matching constraint on $\prod_{i \in I} \mathcal{O}_{\text{life}}^{(i)}$ then each local GFlowNet satisfies the strong flow-matching constraint on $\mathcal{O}_{\text{life}}^{(i)}$.

• if the multi-agent environment is a game environment with Asynchronous Unanimous STOP and if each local GFlowNets satisfy the strong flow-matching constraint on $\mathcal{O}_{life}^{(i)}$ then \hat{R} = $\prod_{i \in I} \hat{R}^{(i)}.$

1120 **Proof 3** We simply define $\mathbb{F} = (\pi^*, F_{\text{out}}^*, F_{\text{init}})$ by $\pi^*(s) \coloneqq (\prod_{i \in I} \pi^{(i),*}(o^{(i)})) / \sim$ ie the projection on \mathcal{A} of the policy toward $\prod_{i \in I} \mathcal{A}^{(i)}$, then F_{out}^* as the product of the measures $F_{\text{out}}^{(i),*}$. Then we define $\tilde{T} = \prod_{i \in I} T^{(i)}$ so that $F_{\text{in}}^*(\prod_{i \in I} \mathcal{A}^{(i)}) = \prod_{i \in I} F_{\text{in}}^{(i),*}(\mathcal{A}^{(i)})$ and $F_{\text{init}} \coloneqq \prod_{i \in I} (F_{\text{in}}^{(i),*} + F_{\text{init}}^{(i)}) - \prod_{i \in I} F_{\text{in}}^{(i),*}$ as the product measure of the $F_{\text{init}}^{(i)}$. By construction this GFlowNet is star-split. 1121 1122 1123 1124

1125 Assume that \mathbb{F} satisfies the strong flow-matching constraint. It follows that for any $A^{(i)} \subset \mathcal{O}_{life}^{(i)}$ we 1126 have 1127 (i)).

$$\prod_{i \in I} F_{\text{in}}^{(i)}(A^{(i)}) = \prod_{i \in I} F_{\text{out}}^{(i)}(A^{(i)}) = \prod_{i \in I} F_{\text{out}}^{(i),*}(A^{(i)})$$

1129 Since, by assumption, all local GFlowNets satisfy the weak flow-matching constraint, all terms in 1130 the left-hand side product are bigger than those in the right-hand side product. Equality may only occur if some term is zero on both sides or if for all $i \in I$, $F_{in}^{(i)} = F_{out}^{(i)}$. Since we assume that the $F_{out}^{(i),*} \neq 0$ we may take all the $A^{(i)} = \mathcal{O}_{life}^{(i)}$ except one to ensure we are in the later case. We conclude 1131 1132 1133 that the strong flow-matching constraint is satisfied for all local GFlowNets on $\mathcal{O}_{ ext{life}}^{(i)}$

1134 If the strong flow-matching constraint is satisfied on $\mathcal{O}_{life}^{(i)}$, then $\hat{R}^{(i)} = R^{(i)} = 0$ on $\mathcal{O}_{life}^{(i)}$. By 1135 construction, $F_{out}^{(i),*} = F_{init}^{(i),*} = 0$ on $\mathcal{O}_{purgatory}^{(i)}$. Therefore, on purgatory, we have

$$\hat{R} = F_{\text{in}} - F_{\text{out}} = F_{\text{in}}^* - F_{\text{out}}^* = \prod_{i \in I} F_{\text{in}}^{(i),*} - \prod_{i \in I} F_{\text{out}}^{(i),*} = \prod_{i \in I} F_{\text{in}}^{(i),*} = \prod_{i \in I} \hat{R}^{(i)}$$

1140 B ALGORITHMS

1137 1138 1139

1141

1144 1145

1177 1178

1180

1187

Algorithm 3 shows the training phase of the independent flow network (IFN). In the each round of IFN, the agents first sample trajectories with policy

$$o_t^{(i)} = p_i(s_t^{(i)}) \text{ and } \pi^{(i)}(o_t^{(i)} \to a_t^{(i)}), \ i \in I$$
 (28)

with $a_t = (a_t^{(i)} : i \in I)$ and $s_{t+1} = T(s_t, a_t)$. Then we train the sampling policy by minimizing the FM loss $\mathcal{L}_{FM}^{\text{stable}}(\mathbb{F}^{(i),\theta})$ for $i \in I$.

	FMI	
Algor	thm 3 Independent Flow Network Training Algorithm for MA-GFlowNets	_
Inpu	: Number of agents N, A multi-agent environment $(S, A, O^{(i)}, A^{(i)}, p_i, S, T, R)$.	
Inpu	: Local GFlowNets $(\pi^{(i),*}, F^{(i),*}_{out}, F^{(i)}_{ioit})_{i \in I}$ parameterized by θ .	
whi	le not converged do	
	Sample and add trajectories $(s_t)_{t\geq 0} \in \mathcal{T}$ to replay buffer with policy according to equation	28
	Generate training distribution of observations $\nu_{\text{state}}^{(i)}$ for $i \in I$ from train buffer	
end	Apply minimization step of FM-loss $\mathcal{L}_{FM}^{\text{stable}}(F_{\text{action}}^{(i),\theta}, R^{(i)})$ for $i \in I$. while	
Algor of CJI	thm 4 shows the training phase of Conditioned Joint Flow Network (CJFN). In the each rou N, we first sample sample trajectories with policy	nd
	$o_t^{(i)} = p_i(s_t^{(i)}) \text{ and } \pi_{\omega}^{(i)}(o_t^{(i)} \to a_t^{(i)}), \ i \in I$ (2)	29)
with a	$t_{i} = (a_{i}^{(i)} : i \in I)$ and $s_{t+1} = T(s_t, a_t)$. Then we train the sampling policy by minimizing t	the
FM lo	ss $\mathbb{E}_{\omega} \mathcal{L}_{\text{FM}}^{\text{stable}}(F_{\text{action}}^{\theta,\text{joint}}(\cdot;\omega), R).$	
Algor	thm 4 Conditioned Joint Flow Network Training Algorithm for MA-GFlowNets	_
Inpu Inpu	: Number of agents N, A multi-agent environment $(S, A, O^{(i)}, A^{(i)}, p_i, S, T, R)$. : Simple Random distribution (Ω, \mathbb{P})	
Inpu	: Local GFlowNets $(\pi^{(i),*}, F^{(i),*}_{out}, F^{(i)}_{ioit})_{i \in I}$ parameterized by θ and $\omega \in \Omega$.	
whi	le not converged do	
	Sample $\omega_1, \dots, \omega_b \sim \mathbb{P}$ and then trajectories $(s_t^{\omega})_{t \geq 0} \in \mathcal{T}$ to replay buffer with policy according	ng
to e	$\{uation 29 \text{ for } \omega \in \{\omega_1, \dots, \omega_b\}$	
	Generate training distribution of states/omega $\nu_{\text{state}}^{\text{M}}$ from the train buffer	
~	Apply minimization step of the FM loss $\mathbb{E}_{\omega} \mathcal{L}_{FM}^{\text{stable}}(\mathbb{F}^{\Theta,\text{joint}}(\cdot;\omega))$ under the constraint of We	ak
flov	-matching.	
end	while	

1179 C DISCUSSION: RELATIONSHIP WITH MARL

Interestingly, there are similar independent execution algorithms in the multi-agent reinforcement learning scheme. Therefore, in this subsection, we discuss the relationship between flow conservation networks and multi-agent RL. The value decomposition approach has been widely used in multi-agent RL based on IGM conditions, such as VDN and QMIX. For a given global state *s* and joint action *a*, the IGM condition asserts the consistency between joint and local greedy action selections in the joint action-value $Q_{\text{tot}}(s, a)$ and individual action values $[Q_i(o_i, a_i)]_{i=1}^k$:

$$\arg\max_{a\in\mathcal{A}}Q_{\text{tot}}(s,a) = \left(\arg\max_{a_1\in\mathcal{A}_1}Q_1(o_1,a_1),\cdots,\arg\max_{a_k\in\mathcal{A}_k}Q_k(o_k,a_k)\right), \forall s\in\mathcal{S}.$$
 (30)

Assumption 1 For any complete trajectory in an MADAG $\tau = (s_0, ..., s_f)$, we assume that $Q_{tot}^{\mu}(s_{f-1}, a) = R(s_f)f(s_{f-1})$ with $f(s_n) = \prod_{t=0}^n \frac{1}{\mu(a|s_t)}$.

Remark 1 Although Assumption 1 is a strong assumption that does not always hold in practical environments. Here we only use this assumption for discussion analysis, which does not affect the performance of the proposed algorithms. A scenario where the assumption directly holds is that we sample actions according to a uniform distribution in a tree structure, i.e., $\mu(a|s) = 1/|\mathcal{A}(s)|$. The uniform policy is also used as an assumption in Bengio et al. (2021).

Lemma 3 Suppose Assumption 1 holds and the environment has a tree structure, based on Theorem 2 and IGM conditions we have:

1199 1) $Q_{tot}^{\mu}(s,a) = F(s,a)f(s);$

1196

1200 2) $(\arg \max_{a_i} Q_i(o_i, a_i))_{i=1}^k = (\arg \max_{a_i} F_i(o_i, a_i))_{i=1}^k$.

Based on Assumption1, we have Lemma 3, which shows the connection between Theorem 2 and the IGM condition. This action-value function equivalence property helps us better understand the multi-flow network algorithms, especially showing the rationality of Theorem 2.

1205 1206 С.1 Ркооf of Lemma 3

Proof 4 The proof is an extension of that of Proposition 4 in Bengio et al. (2021). For any (s, a)satisfies $s_f = T(s, a)$, we have $Q_{tot}^{\mu}(s, a) = R(s_f)f(s)$ and $F(s, a) = R(s_f)$. Therefore, we have $Q_{tot}^{\mu}(s, a) = F(s, a)f(s)$. Then, for each non-final node s', based on the action-value function in terms of the action-value at the next step, we have by induction:

1216 1217 where $\hat{R}(s')$ is the reward of $Q_{tot}^{\mu}(s,a)$ and (a) is due to that $\hat{R}(s') = 0$ if s' is not a final state. 1218 Since the environment has a tree structure, we have

 $F(s,a) = \sum_{a' \in \mathcal{A}(s')} F(s',a'), \tag{32}$

1221 which yields

$$Q^{\mu}_{tot}(s,a) = \mu(a|s')F(s,a)f(s') = \mu(a|s')F(s,a)f(s)\frac{1}{\mu(a|s')} = F(s,a)f(s)$$

1225 According to Theorem 2, we have $F(s_t, a_t) = \prod_i F_i(o_t^i, a_t^i)$, yielding

$$\arg\max_{a} Q_{tot}(s,a) \stackrel{(a)}{=} \arg\max_{a} \log F(s,a)f(s)$$

$$\stackrel{(b)}{=} \arg\max_{a} \sum_{i=1}^{k} \log F_{i}(o_{i},a_{i})$$

$$\stackrel{(c)}{=} \left(\arg\max_{a_{1}\in\mathcal{A}_{i}} F_{1}(o_{1},a_{1}), \cdots, \arg\max_{a_{k}\in\mathcal{A}_{k}} F_{k}(o_{k},a_{k})\right),$$
(33)

where (a) is based on the fact F and f(s) are positive, (b) is due to Theorem 2. Combining with the IGM condition

$$\arg\max_{a\in\mathcal{A}}Q_{tot}(s,a) = \left(\arg\max_{a_1\in\mathcal{A}_1}Q_1(o_1,a_1),\cdots,\arg\max_{a_k\in\mathcal{A}_k}Q_k(o_k,a_k)\right), \forall s\in\mathcal{S}.$$
 (34)

1238 we can conclude that

1239
1240
$$\left(\arg\max_{a_i\in\mathcal{A}_i}F_i(o_i,a_i)\right)_{i=1}^k = \left(\arg\max_{a_i\in\mathcal{A}_1}Q_i(o_i,a_i)\right)_{i=1}^k$$
1241

Then we complete the proof.

1233 1234 1235

1236 1237

1219 1220

1222 1223 1224

1242 D ADDITIONAL EXPERIMENTS

1244 D.1 Hyper-Grid Environment

1246 D.1.1 EFFECT OF SAMPLING METHOD:

1248 We consider two different sampling methods of JFN; the first one is to sample trajectories using the 1249 flow function F_i of each agent independently, called JFN (IS), and the other one is to combine the 1250 policies π_i of all agents to obtain a joint policy π , and then performed centralized sampling, named 1251 JFN (CS). As shown in Figure 7, we found that the JFN (CS) method has better performance than 1252 JFN (IS) because the error of the policy π estimated by the combination method is smaller, and 1253 several better samples can be obtained during the training process. However, the JFN (IS) method 1254 can achieve decentralized sampling, which is more in line with practical applications.



Figure 7: The performance of JFN with different methods.

1271 D.1.2 EFFECT OF DIFFERENT REWARDS:

We study the effect of different rewards in Figure 8. In particular, we set $R_0 = \{10^{-1}, 10^{-2}, 10^{-4}\}$ for different task challenge. A smaller value of R_0 makes the reward function distribution more sparse, which makes policy optimization more difficult Bengio et al. (2021); Riedmiller et al. (2018); Trott et al. (2019). As shown in Figure 8, we found that our proposed method is robust with the cases $R_0 = 10^{-1}$ and $R_0 = 10^{-2}$. When the reward distribution becomes sparse, the performance of the proposed algorithm degrades slightly.



Figure 8: The effect of different reward R_0 on different algorithm.

1293 D.1.3 FLOW MATCH LOSS FUNCTION:



For some RL algorithms based on the state-action value function estimation, the loss usually oscillates. This may be because RL-based methods use experience replay buffer and the transition data distribution is not stable enough. The method we propose uses an on-policy based optimization method, and the data distribution changes with the current sampling policy, hence the loss function is relatively stable. Then we present the experimental details on the Hyper-Grid environments. We set the same number of training steps for all algorithms for a fair comparison. Moreover, we list the key hyperparameters of the different algorithms in Tables 2-6.



Figure 9: The flow matching loss of different algorithm.

In addition, as shown in Table 1, both the reinforcement learning methods and our proposed method can achieve the highest reward, but the average reward of reinforcement learning is slightly better for all found modes. Our algorithms do not always have higher rewards compared to RL, which is reasonable since the goal of MA-GFlowNets is not to maximize rewards.

Environment	MAPPO	MASAC	MCMC	CFN	JFN
Hyper-Grid v1	2.0	1.84	1.78	2.0	2.0
Hyper-Grid v2	1.90	1.76	1.70	1.85	1.85
Hyper-Grid v3	1.84	1.66	1.62	1.82	1.82

Table 1: The best reward found using different methods.

1334 D.2 STARCRAFT

1303 1304

1305

1306

1309 1310

1311

1313

1315

1316

1317 1318

1323 1324

1326

1328

1335

We present a visual analysis based on 3m with three identical entities attacking to win. All compari-1336 son experiments adopted PyMARL framework and used default experimental parameters. Figure 10 1337 shows the decision results of different algorithms on the 3m map. It can be found that the proposed 1338 algorithm can obtain results under different reward distributions, that is, win at different costs. The 1339 costs of other algorithms are often the same, which shows that the proposed algorithm is suitable 1340 for scenarios with richer rewards. Figure 11 shows the performance of the different algorithms on 1341 2s3z, which shows a similar conclusion that the algorithm based on GFlowNets may be difficult 1342 to get the best yield, but the goal is not to do this, but to fit the distribution better. Moreover, on 1343 StarCraft missions, we did not use a clear metric to indicate the diversity of different trajectories, 1344 mainly because the status of each entity includes multiple aspects, its movement range, health, op-1345 ponent observation, etc., which can easily result in different trajectories, but these differences do not indicate a good fit for the reward distribution. As a result, it is not presented in the same way as Hyper-Grid and Simple-Spread. Therefore, we used a visual method to compare the results. The 1347 maximized reward-oriented algorithms such as QMIX will improve the reward by reducing the death 1348 of entities, while the GFlowNets method can better fit the distribution on the basis of guaranteeing 1349 higher rewards.



1403 In each better.



Figure 12: Average return and the number of distinctive trajectories performance of different algorithms on Sparse-Simple-Spread environments.

Table 2: Hyper-parameter of MAPPO under different environments

	Hyper-Grid-v1	Hyper-Grid-v2	Hyper-Grid-v3
Train Steps	20000	20000	20000
Agent	2	2	3
Grid Dim	2	3	3
Grid Size	[8,8]	[8,8]	[8,8]
Actor Network Hidden Layers	[256,256]	[256,256]	[256,256]
Optimizer	Adam	Adam	Adam
Learning Rate	0.0001	0.0001	0.0001
Batchsize	64	64	64
Discount Factor	0.99	0.99	0.99
PPO Entropy	1e-1	1e-1	1e-1

Table 3: Hyper-parameter of MASAC under different environments

	Hyper-Grid-v1	Hyper-Grid-v2	Hyper-Grid-v3
Train Steps	20000	20000	20000
Grid Dim	2	3	3
Grid Size	[8,8]	[8,8]	[8,8]
Actor Network Hidden Layers	[256,256]	[256,256]	[256,256]
Critic Network Hidden Layers	[256,256]	[256,256]	[256,256]
Optimizer	Adam	Adam	Adam
Learning Rate	0.0001	0.0001	0.0001
Batchsize	64	64	64
Discount Factor	0.99	0.99	0.99
SAC Alpha	0.98	0.98	0.98
Target Network Update	0.001	0.001	0.001

Table 4: Hyper-parameter of JFN under different environments

	Hyper-Grid-v1	Hyper-Grid-v2	Hyper-Grid-v3
Train Steps	20000	20000	20000
R_2	2	2	2
R_1	0.5	0.5	0.5
Grid Dim	2	3	3
Grid Size	[8,8]	[8,8]	[8,8]
Trajectories per steps	16	16	16
Flow Network Hidden Layers	[256,256]	[256,256]	[256,256]
Optimizer	Adam	Adam	Adam
Learning Rate	0.0001	0.0001	0.0001
ϵ	0.0005	0.0005	0.0005

Table 5: Hyper-parameter of CJFN under different environments				
	Hyper-Grid-v1	Hyper-Grid-v2	Hyper-Grid-v3	
Train Steps	20000	20000	20000	
R_2	2	2	2	
R_1	0.5	0.5	0.5	
Grid Dim	2	3	3	
Grid Size	[8,8]	[8,8]	[8,8]	
Trajectories per steps	16	16	16	
Flow Network Hidden Layers	[256,256]	[256,256]	[256,256]	
Optimizer	Adam	Adam	Adam	
Learning Rate	0.0001	0.0001	0.0001	
ϵ	0.0005	0.0005	0.0005	
Number of ω	4	4	4	

Table 6: Hyper-parameter of CFN under different environments

	Hyper-Grid-v1	Hyper-Grid-v2	Hyper-Grid-v3
Train Steps	20000	20000	20000
Trajectories per steps	16	16	16
R_2	2	2	2
R_1	0.5	0.5	0.5
Grid Dim	2	3	3
Grid Size	[8,8]	[8,8]	[8,8]
Flow Network Hidden Layers	[256,256]	[256,256]	[256,256]
Optimizer	Adam	Adam	Adam
Learning Rate	0.0001	0.0001	0.0001
ϵ	0.0005	0.0005	0.0005