



HAL
open science

Persuasion with limited communication capacity

Mael Le Treust, Tristan Tomala

► **To cite this version:**

Mael Le Treust, Tristan Tomala. Persuasion with limited communication capacity. Journal of Economic Theory, 2019. hal-01633656v3

HAL Id: hal-01633656

<https://hal.science/hal-01633656v3>

Submitted on 4 Sep 2019

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Persuasion with limited communication capacity*

Maël Le Treust[†] and Tristan Tomala[‡]

Accepted in JET, September 1, 2019

Abstract

We consider a Bayesian persuasion problem where the persuader and the decision maker communicate through an imperfect channel that has a fixed and limited number of messages and is subject to exogenous noise. We provide an upper bound on the payoffs the persuader can secure by communicating through the channel. We also show that the bound is tight, i.e., if the persuasion problem consists of a large number of independent copies of the same base problem, then the persuader can achieve this bound arbitrarily closely by using strategies that tie all the problems together. We characterize this optimal payoff as a function of the information-theoretic capacity of the communication channel.

Keywords: Bayesian persuasion, communication channel, mutual information.

JEL Classification Numbers: C72, D82, D83.

*The authors thank James Best, Olivier Gossner, Frédéric Koessler, Marie Laclau, Daniel Martin, Ludovic Renou, Thomas Rivera, Jakub Steiner, Colin Stewart for stimulating discussions and comments. We thank the editor Alessandro Pavan, an anonymous associate editor and anonymous referees for helpful comments and suggestions. We also thank participants of the 6th workshop on Stochastic Methods in Game Theory, Erice May 2017; the 13th European Meeting on Game Theory (SING13), Paris July 2017; the XXVI Colloque Gretsi, Juan-Les-Pins, September 2017; the 10th Transatlantic Theory Workshop, Paris September 2017; the 55th Allerton Conference, Monticello, Illinois, October 2017. We thank the Institute Henri Poincaré for hosting numerous research meetings.

[†]ETIS UMR 8051, Université Paris Seine, Université Cergy-Pontoise, ENSEA, CNRS, F-95000, Cergy, France; mael.le-treust@ensea.fr; sites.google.com/site/maelletreust/. This research has been conducted as part of the project Labex MME-DII (ANR11-LBX-0023-01). Maël Le Treust gratefully acknowledges financial support from INS2I CNRS, DIM-RFSI, SRV ENSEA, The Paris Seine Initiative and IEA Cergy-Pontoise.

[‡]HEC Paris and GREGHEC, 1 rue de la Libération, 78351 Jouy-en-Josas, France; tomala@hec.fr; sites.google.com/site/tristantomala2. Tristan Tomala gratefully acknowledges the support the HEC foundation and ANR/Investissements d'Avenir under grant ANR-11-IDEX-0003/Labex Ecodec/ANR-11-LABX-0047.

1 Introduction

In modern internet societies, pieces of information are repeatedly and continuously disclosed to decision makers by informed agents. Information transmission is affected by at least two sources of friction. First, the sender and the receiver of a given message may have nonaligned incentives, in which case the sender might be unwilling to transmit truthful information. Second, communication between agents is often imperfect. The sender and the receiver may have time constraints to write or read messages, forcing the sender to summarize his arguments and making him unable to convey all the details. Further, there might be discrepancies between the informational content of a message that is intended by the sender and the one understood by the receiver. For instance, if the mother tongue of the sender and of the receiver are different, there are possible translation errors (See [Blume, Board, and Kawamura, 2007](#)). Additionally, messages travelling in a network of computers might be subject to random shocks, internal errors or protocol failures. Studying the effect of noise in communication channels is the starting point of information theory ([Shannon, 1948](#)).

Our paper aims to study the following questions. How does imperfect communication reduce the possibilities of persuasion in a sender-receiver interaction? When the sender communicates many pieces of information, to what extent does tying the pieces together help in overcoming the communication limitations?

We consider a sender and a receiver who communicate over an imperfect channel and are engaged in a series of $n \geq 1$ persuasion problems. The sender observes n independent and identically distributed pieces of information and sends $k \geq 1$ messages to the receiver. Messages are sent through a channel that consists of two finite sets X, Y of respectively inputs and outputs messages and of a transition probability Q from X to Y such that when the sender chooses input message x , the receiver receives output message y with probability $Q(y|x)$. Upon receiving k output messages from the channel, the receiver chooses n actions, one for each problem. Payoffs are additively separable across persuasion problems. We assume that the sender is able to commit to a disclosure strategy that maps sequences of pieces of information to distributions of sequences of input messages.

We study the optimal average payoff secured by the sender by committing to a strategy. We give an upper bound on this optimal payoff and show that this bound is achieved asymptotically

when the numbers n and k grow large. To prove this latter statement, we borrow techniques from information theory, namely, the coding and decoding schemes of [Shannon \(1948, 1959\)](#). This machinery allows to transmit a sequence of messages over a noisy channel with the property that the receiver recovers almost all messages correctly. The information theoretic literature typically considers an obedient receiver who calculates the decoded messages and takes them at face value. In the persuasion game framework, the receiver is strategic and may not follow any prescribed scheme. Rather, the receiver takes into account the strategy of the sender and the received outputs, calculates its Bayesian belief about the sequence of states, and chooses a sequence of actions that maximizes its payoff. Our technical contribution is to construct a strategy of the sender for which we are able to estimate and to control those Bayesian beliefs in order to ensure that the strategic receiver chooses a desired sequence of actions.

Our upper bound is the value of an *optimal splitting problem with information constraint*, which represents the best payoff that the sender can achieve by sending a message, subject to the constraint that the mutual information between the state and the message is no more than the capacity of the channel. We show that this value is given by the concave closure of the payoff function of the sender, subject to a constraint on the entropy of posterior beliefs. This is also given by the concave closure of a modified payoff function, where the sender pays a cost proportional to the mutual information between the state and the message.

1.1 Motivating example.

There are relevant situations where a sender discloses information about a large number of independent state parameters. For instance, one can think of testing product quality: a firm has many items to sell, which are ex-ante identical, and the authorities (e.g., the FDA for drugs) design quality tests ¹. One can also think about designing and grading exams to assess the quality of a large number of students².

As an example, consider an innovating firm that has several projects to be financed by investors. The board of investors audits the firm, which is given a limited amount of time to present all the projects. How to best structure arguments in order to get the maximum number of projects approved?

To be specific, let us assume that all projects are ex-ante identical and equally likely to be

¹See e.g., [Perez and Skreta, 2018](#).

²See [Boleslavsky and Cotton, 2015](#) for a model of grading standards through Bayesian persuasion.

of good or bad quality. When a project is approved, it yields a positive return of $+1$ to the investors if it is good, and a negative return of -7 if it is bad; rejecting a project yields a payoff of 0 . The objective of the firm is to get a maximum number of projects approved.

Suppose that the firm commits to an information disclosure mechanism, i.e., distributions of messages conditional on states (as in [Kamenica and Gentzkow, 2011](#)) and faces no restriction on the number of messages. To invest, the board of investors must be persuaded that the project is good with probability at least $7/8$. Thus, for each project, the firm would optimally draw a good message g or a bad message b with the following probabilities:

$$\mathbb{P}(g \mid \text{project is good}) = 1, \quad \mathbb{P}(g \mid \text{project is bad}) = 1/7.$$

This way, the belief that the project is good upon receiving the good message is as follows:

$$\mathbb{P}(\text{project is good} \mid g) = 7/8,$$

and the project is accepted with probability $4/7$ (see [Section 4](#)).

Now, suppose that the auditing board gives the firm only half the time it would require to talk about all projects. Namely, there is an even number n of projects, but the firm has only $n/2$ messages available.

A simple strategy the firm can adopt would be to select half of the projects, focus on them, and communicate optimally for each of them. With this strategy, half of the projects are accepted with probability $4/7$ each, so in expectation, the average number of accepted projects is $2/7$. This is not optimal, and a better strategy would be to pair projects by two and to draw one message g, b for each pair in the following way:

$$\mathbb{P}(g \mid \text{both projects are good}) = 1, \quad \mathbb{P}(g \mid \text{both projects are bad}) = 0,$$

$$\mathbb{P}(g \mid \text{only one project is good}) = 1/6.$$

The total probability of g is $1/3$ and upon observing this message, the beliefs about quality are as follows:

$$\mathbb{P}(\text{both projects are good} \mid g) = 6/8,$$

$$\mathbb{P}(\text{only project 1 is good} \mid g) = \mathbb{P}(\text{only project 2 is good} \mid g) = 1/8.$$

Therefore, each project is believed to be good with probability $7/8$ and both projects are accepted when g is received. Thus, the expected average number of accepted projects is $1/3 > 2/7$.

We thus see that tying projects together improves upon communication about each project separately. Suppose that the number of projects is large. Is it possible to find a more complex strategy that further improves the payoff?

Our main result, Theorem 3.1, gives an upper bound on the expected average number of accepted projects when the number of messages is half the number of projects. The upper bound is tight: the optimal value approaches it as the number of project increases. In this example, the upper bound is λ^* where (λ^*, p^*) is the unique solution in $[0, 1] \times [0, \frac{1}{2}]$ of the system of equations:

$$\frac{1}{2} = \lambda^* \frac{7}{8} + (1 - \lambda^*)p^*, \quad \frac{1}{2} = \lambda^* H\left(\frac{7}{8}\right) + (1 - \lambda^*)H(p^*),$$

where $H(p) = -p \log(p) - (1 - p) \log(1 - p)$ is the entropy function. The first equation is Bayes plausibility (Kamenica and Gentzkow, 2011) coming from Bayes' rule, saying that the expected posterior belief is the prior belief. The second equation requires the expected entropy of the posterior to be $\frac{1}{2}$, which means that the *mutual information* between the quality of the project and the message sent to the receiver is equal to the number of messages per project that the firm is able to transmit.

Numerically $\lambda^* \approx 0.519 < \frac{4}{7} \approx 0.571$. Thus, for large n , the sender can achieve a payoff better than $1/3$ but bounded away from the payoff obtained with unrestricted communication.

1.2 Related literature

We now describe the relationships between our contribution and the literature. This paper is at the junction of Bayesian persuasion and information theory.

The traditional game theoretic approach to strategic information disclosure assumes perfect communication and analyzes in isolation the problem of sending a single message. These are the well-known sender-receiver games where an informed player, the sender, communicates once with a receiver who takes an action. In the *cheap talk* version of this game, the message sent by the sender is costless and unverifiable; see for instance the seminal paper of Crawford and Sobel (1982). In the *Bayesian persuasion* game (Kamenica and Gentzkow, 2011), the sender chooses verifiably an information disclosure device prior to learning his information. That is,

the sender is an *information designer* (Bergemann and Morris, 2016, 2017; Taneva, 2018) who chooses, without knowledge of the state, the information or signaling structure which releases information to the decision maker.

In parallel, information theory considers agents with perfectly aligned interests and analyzes the *rate* of information transmission. The sender observes an information *flow*, which is a stochastic process, and sends messages to the receiver over an imperfect channel represented by a transition probability from input to output messages. Truthful information transmission is the common goal of the sender and the receiver. The *rate of information transmission* is the average number of correct guesses made by the receiver. Shannon's theory (Shannon, 1948, 1959) determines whether a source of information can be transmitted over the channel with arbitrarily small probability of error and shows that the rate of the source of information has to be smaller than *the capacity of the channel* defined as the maximal *mutual information* between input and output messages.

Our model of persuasion has two essential features. The sender and the receiver are engaged in a large number of identical copies of the same game and communication is restricted to an imperfect channel. As Kamenica and Gentzkow (2011), we consider the payoff obtained by the sender as a function of the belief of the receiver, when the receiver takes optimal actions. With unrestricted communication, that is on a perfect channel with large set of inputs, the optimal payoff for the sender is given by the concave closure of this function. Then, solving any number of identical games amounts to solving each copy separately. With a single copy, the game of persuasion with a noisy channel is studied by Tsakas and Tsakas (2018) who prove the existence of optimal solutions and show monotonicity of the sender's payoff with respect to the noise of the channel. Considering many copies of the base game *and* restricted communication, we show that *linking independent problems together* yields a better payoff to the sender: the optimal strategy correlates all messages with the state parameters of all problems. In this respect, our work bears some similarity with Jackson and Sonnenschein (2007), who showed that a mechanism designer can achieve more outcomes in an incentive compatible manner by linking many identical problems together.

The optimal payoff that we characterize is related to models where the cost of information is measured by mutual information. Such information costs have been introduced in the literature

on rational inattention by [Sims \(2003\)](#), (See also [Martin, 2017](#); [Matejka and McKay, 2015](#); [Steiner, Stewart, and Matejka, 2017](#)). The use of mutual information has been axiomatized in [Morris and Strack \(2019\)](#) and [Hebert and Woodford \(2018\)](#). In the context of persuasion, [Gentzkow and Kamenica \(2014\)](#) consider a model where the sender gets his payoff from the game, minus a cost that is proportional to the mutual information between the state and the message; see also [Matyskova \(2018\)](#). With Lagrangian methods, we find that the value of our optimal splitting problem with information constraint is the concave closure of the payoff function, net of such an information cost, a similar concavification problem is found in [Caplin and Dean \(2013\)](#)

Different from those papers, the mutual information is not a primitive of our model. Our finding is that the noise and limitations in communication induce a *shadow cost* measured by the mutual information.

Entropy and mutual information appear endogenously in several papers on repeated games where players have bounded rationality ([Neyman and Okada, 1999, 2000](#)), are not able to freely randomize their actions ([Gossner and Vieille, 2002](#)), or observe actions imperfectly ([Gossner and Tomala, 2006, 2007](#)). A related paper is [Gossner, Hernández, and Neyman \(2006\)](#), henceforth GHN, who also consider a sender-receiver game. In GHN, the sender and the receiver play an infinitely repeated game with common interests: both the sender and the receiver want to choose the action that matches the state. The sender knows the infinite sequence of states and can communicate with the receiver only through his actions. GHN characterize the best average payoff that the sender (and the receiver) can achieve. Their solution resembles ours: the optimal value is the payoff obtained when the sender can send a direct message to the receiver, subject to an information constraint.

There are important differences with our work. First, GHN study a cheap talk game with common interests. By contrast, we do not assume common interests and we assume commitment power for the sender. Second, GHN is a truly repeated game model: at any given time t , both players choose actions and the information of the receiver at this time consists of past actions. In our case, the sender knows a finite sequence of states and chooses a finite sequence of input messages, the receiver observes a finite sequence of output messages and chooses a sequence of actions. This is why, rather than seeing our model as a repeated game of persuasion, we view it as a *spatial* model with identical copies of the same problem coexisting at the same time. This also explains why the number of copies n need not be equal to the number of times k the channel

is used by the sender. Our result characterizes the optimal payoff as a function of the ratio of the number n of pieces of information to the number k of channel uses. In particular, this allows us to analyze cases where the channel is perfect (i.e. not subject to random noise) but with limited input size: there are fewer messages than states or actions.

Cheap talk with a noisy channel has been studied by [Blume, Board, and Kawamura \(2007\)](#) who show that the presence of noise is possibly welfare improving. Such a phenomenon cannot happen in the persuasion context as the sender could commit to replicate the noise. Relatedly, [Hernández and von Stengel \(2014\)](#) consider a sender-receiver game with common interests over an imperfect channel. In that paper, there is only one state known by the sender and one action taken by the receiver, while the channel can be used a fixed number of times. [Hernández and von Stengel \(2014\)](#) characterize all the Nash equilibria of this game and study the differences with Shannon’s coding methods. Again, we do not assume common interests and assume commitment power for the sender. More importantly, our focus is different and more in line with GHN: we do not treat a single persuasion problem but a large sequence of them and use information theory to study the asymptotics of the problem.

Our work is also related to some information theoretic literature. Following GHN, a line of papers study *empirical coordination* between a sender and a receiver ([Cuff, Permuter, and Cover, 2010](#); [Cuff and Zhao, 2011](#); [Le Treust, 2017](#)). Assuming common interest between the sender and the receiver, those papers characterize the asymptotic empirical distributions of (states, messages, actions) which are achievable, given the information structure and the noisy channel. The closest paper in this literature is [Le Treust and Tomala \(2016\)](#) where we have studied empirical coordination between a persuader and a decision maker induced by approximate equilibria as the number of repetitions tends to infinity. Recently, [Akyol, Langbort, and Başar \(2017\)](#) have considered the problem of Bayesian persuasion in a model with Gaussian states and channel and quadratic functions as in [Crawford and Sobel \(1982\)](#).

The remainder of this paper is organized as follows. The model is described in [Section 2](#) and we state our main results in [Section 3](#). In [Section 4](#), we illustrate our results with a detailed example. We provide an extension in [Section 5](#) and concluding comments in [Section 6](#). Proofs are in the Appendix.

2 Model

2.1 The persuasion problem

In this model, we consider a sender (S) and a receiver (R) engaged in a series of identical persuasion problems and where the communication technology is fixed exogenously.

There is a finite state space Ω endowed with a prior probability distribution μ , a finite action set A for the receiver, and each player $i = S, R$ has a payoff function $u_i : \Omega \times A \rightarrow \mathbb{R}$. There is also a fixed communication channel (X, Y, Q) , where X, Y are finite sets of messages and $Q : X \rightarrow \Delta(Y)$ is a transition probability from X to Y (henceforth $\Delta(S)$ denotes the set of probability distributions over the finite set S).

Given two integers n, k , we define a repeated persuasion problem where the uncertainty is about a sequence $\omega^n = (\omega_1, \dots, \omega_n)$ drawn i.i.d. from (Ω, μ) . The receiver chooses a sequence of actions $a^n = (a_1, \dots, a_n)$ and the payoff for player $i = S, R$ is as follows:

$$\bar{u}_i(\omega^n, a^n) = \frac{1}{n} \sum_{t=1}^n u_i(\omega_t, a_t).$$

To disclose information, the sender can use the channel k times by choosing a sequence of input messages $x^k = (x_1, \dots, x_k)$. The channel then draws a sequence of output messages y^k with probability $Q^k(y^k|x^k) = \prod_{t=1}^k Q(y_t|x_t)$ and sends it to the receiver.

This defines the following persuasion game $\Gamma(n, k)$:

1. The sender chooses a strategy $\sigma : \Omega^n \rightarrow \Delta(X^k)$ which is announced to the receiver.
2. A sequence of states ω^n is drawn i.i.d. from the prior μ , a sequence of input messages x^k is drawn with probability $\sigma(x^k|\omega^n)$, a sequence of output messages y^k is drawn with probability $Q^k(y^k|x^k)$ and is observed by the receiver.
3. The receiver chooses a sequence of actions a^n .

Then, player $i = S, R$ gets the average payoff $\bar{u}_i(\omega^n, a^n)$.

Notice that for $n = k = 1$, this is the model of [Tsakas and Tsakas \(2018\)](#) of a single persuasion problem with noisy communication. An interesting particular case is given by *perfect* channels where $X = Y$ and $Q(y|x) = \mathbb{1}_{\{y=x\}}$. In such a case, the only limitation is given by the number of messages. If we let $n = k = 1$ and choose a perfect channel with sufficiently many messages $|X| =$

$|Y| \geq |\Omega|$, the model encompasses the standard persuasion game of [Kamenica and Gentzkow \(2011\)](#).

2.2 Optimal robust payoff

As a solution concept, we study the best payoff the sender can secure, regardless of which best reply is chosen by the receiver. A strategy of the receiver is a mapping $\tau : Y^k \rightarrow A^n$. Knowing σ , the receiver chooses a best reply τ , which maximizes the expected payoff. That is, for each y^k :

$$\tau(y^k) \in \operatorname{argmax}_{a^n \in A^n} \sum_{\omega^n, x^k} \mu^n(\omega^n) \sigma(x^k | \omega^n) Q(y^k | x^k) \bar{u}_R(\omega^n, a^n).$$

Denote $BR(\sigma)$ the set of best replies of the receiver to the strategy σ .

Definition 2.1. *The optimal robust payoff of the sender in this problem is as follows:*

$$U_S^*(\mu^n, Q^k) = \sup_{\sigma} \min_{\tau \in BR(\sigma)} \sum_{\omega^n, x^k, y^k} \mu^n(\omega^n) \sigma(x^k | \omega^n) Q^k(y^k | x^k) \bar{u}_S(\omega^n, \tau(y^k)).$$

This definition differs from the conventional solution to Bayesian persuasion of [Kamenica and Gentzkow \(2011\)](#) where the receiver takes the best reply which is preferred by the sender. Our choice is motivated by robustness; we ask the solution to be robust to the way the receiver breaks ties³. We stress that this choice does not matter for generic problems. Indeed, with slight perturbations of the payoff function of the receiver, we can make sure that indifferences occur only at interior beliefs. When this is the case, the sender can slightly change his strategy in order to avoid the indifference region.

The goal of this paper is to give an upper bound for the optimal robust payoff and to characterize its limit when n and k tend to infinity.

2.3 Optimal splitting problem with information constraint

To state our main results, we introduce some definitions.

Definition 2.2. *A splitting of $\mu \in \Delta(\Omega)$ is a finite family $(\lambda_m, \nu_m)_m$, where for each m , $\nu_m \in \Delta(\Omega)$, $\lambda_m \in [0, 1]$, $\sum_m \lambda_m = 1$ such that:*

$$\mu = \sum_m \lambda_m \nu_m. \tag{1}$$

³A similar approach is followed by [Inostroza and Pavan \(2018\)](#) and [Mathevet, Perego, and Taneva \(2019\)](#).

A splitting of μ is a distribution of posterior beliefs whose average equals the prior. An “information structure” which draws a message m with probability $\mathbb{P}(m|\omega)$ in state ω , induces a splitting $(\lambda_m, \nu_m)_m$ with $\lambda_m = \sum_{\omega'} \mu(\omega') \mathbb{P}(m|\omega')$ and $\nu_m(\omega) = \frac{\mu(\omega) \mathbb{P}(m|\omega)}{\sum_{\omega'} \mu(\omega') \mathbb{P}(m|\omega')}$. From the splitting lemma (Aumann and Maschler, 1995) or Bayes plausibility (Kamenica and Gentzkow, 2011), for each decomposition of the prior belief into a convex combination of posterior $\mu = \sum_m \lambda_m \nu_m$, the splitting $(\lambda_m, \nu_m)_m$ is induced by some information structure, for example, $\mathbb{P}(m|\omega) = \lambda_m \nu_m(\omega) / \mu(\omega)$.

For each posterior belief $\nu \in \Delta(\Omega)$, let the set of optimal actions of the receiver be:

$$A^*(\nu) = \operatorname{argmax}_{a \in A} \sum_{\omega} \nu(\omega) u_R(\omega, a).$$

We denote by $u_S^*(\nu) = \min_{a \in A^*(\nu)} \sum_{\omega} \nu(\omega) u_S(\omega, a)$ the *robust payoff* of the sender at the belief ν , i.e., the payoff of the sender when the receiver chooses the optimal action, which is *worst* for S .

We now introduce tools borrowed from information theory; the reader is referred to Cover and Thomas (2006).

Definition 2.3. 1. The (Shannon) entropy of a probability distribution $q \in \Delta(S)$ over a finite set S is as follows:

$$H(q) = - \sum_q q(s) \log q(s),$$

where the logarithm has basis 2 and $0 \log 0 = 0$.

2. The mutual information between two random variables (\mathbf{x}, \mathbf{y}) , drawn from the joint probability distribution $p(x)Q(y|x)$ is as follows:

$$I_{p,Q}(\mathbf{x}; \mathbf{y}) = H\left(\sum_x p(x)Q(\cdot|x)\right) - \sum_x p(x)H(Q(\cdot|x))$$

3. The capacity of the channel (X, Y, Q) is as follows:

$$C(Q) = \max_{p \in \Delta(X)} I_{p,Q}(\mathbf{x}; \mathbf{y}).$$

The channel capacity $C(Q)$ is the maximal mutual information between two random variables (\mathbf{x}, \mathbf{y}) , respectively the input and output of the channel, drawn from the joint probability distribution $p(x)Q(y|x)$, where the maximum is over the marginal distribution $p(x)$. Intuitively,

this is the maximal number of bits of information that can be transmitted reliably through the channel (see [Cover and Thomas, 2006](#)).

Equipped with these tools, our main definition is the following.

Definition 2.4. *For any $c \geq 0$, the optimal splitting problem with information constraint is:*

$$\begin{aligned} V(\mu, c) = & \sup \sum_m \lambda_m u_S^*(\nu_m) \\ \text{s.t.} & \sum_m \lambda_m \nu_m = \mu, \\ \text{and} & H(\mu) - \sum_m \lambda_m H(\nu_m) \leq c. \end{aligned}$$

This is the best payoff that the sender can secure by choosing a splitting of the prior belief (i.e., an information structure) under the constraint that the expected reduction of entropy does not exceed the capacity c of the channel. The entropy reduction $H(\mu) - \sum_m \lambda_m H(\nu_m)$ is nonnegative and is the mutual information between a random state ω and a random message \mathbf{m} , drawn from the joint distribution $(\lambda_m \nu_m(\omega))_{(\omega, \mathbf{m})}$. The interpretation is thus that the sender optimizes over a set of information structures that convey bounded information about the state.

Notice that $V(\mu, c)$ is less than or equal to the concave closure (or concavification) of u_S^* at μ which is the unconstrained supremum $\text{cav } u_S^*(\mu) := \sup \{ \sum_m \lambda_m u_S^*(\nu_m) : \sum_m \lambda_m \nu_m = \mu \}$.

3 Results

3.1 The main result

The main result of this paper shows that the value of the optimal splitting problem with information constraint provides an upper bound to the optimal robust payoff and that the bound is achieved asymptotically.

Theorem 3.1. *1. The optimal robust payoff of the sender is no more than the value of the optimal splitting problem with information constraint. For each pair of integers n, k :*

$$U_S^*(\mu^n, Q^k) \leq V\left(\mu, \frac{k}{n} C(Q)\right).$$

2. The optimal robust payoff of the sender converges to the value of the optimal splitting problem with information constraint in the following sense. For each $r \in [0, +\infty]$, for each pair of sequences of integers $(k_j, n_j)_{j \in \mathbb{N}}$ such that $\lim_{j \rightarrow \infty} \max(n_j, k_j) = \infty$ and $\lim_{j \rightarrow \infty} \frac{k_j}{n_j} = r$,

we have:

$$\lim_{j \rightarrow \infty} U_S^*(\mu^{n_j}, Q^{k_j}) = V(\mu, rC(Q)).$$

On the one hand, this result shows communication restrictions limits the payoff that can be achieved through Bayesian persuasion. On the other hand, it quantifies the extent to which repeating the same problem and linking the copies together helps in overcoming those restrictions.

3.1.1 Sketch of proof

We give an intuition for the main arguments of the proof; the technical details are in the appendix.

First point, upper bound. The argument is that regardless of which strategies are used, the mutual information between the states and the messages to the receiver cannot exceed the capacity of the channel.

For simplicity, consider the case $n = k = 1$ where the result says $U_S^*(\mu, Q) \leq V(\mu, C(Q))$. Take any strategy σ of the sender. This induces the splitting $\mu = \sum_y \mathbb{P}_\sigma(y) \nu_y$ where $\mathbb{P}_\sigma(y) = \sum_{\omega, x} \mu(\omega) \sigma(x|\omega) Q(y|x)$ is the probability of the message y and

$$\nu_y(\omega) = \mathbb{P}_\sigma(\omega|y) = \frac{\sum_x \mu(\omega) \sigma(x|\omega) Q(y|x)}{\mathbb{P}_\sigma(y)}$$

is the posterior belief conditional on y . The mutual information of this splitting is:

$$H(\mu) - \sum_y \mathbb{P}_\sigma(y) H(\nu_y) := I(\boldsymbol{\omega}; \mathbf{y})$$

where $(\boldsymbol{\omega}, \mathbf{x}, \mathbf{y})$ denotes a random triple of state, input and output messages drawn from the joint distribution $\mu(\omega) \sigma(x|\omega) Q(y|x)$. With an abuse of notation, we denote $I(\boldsymbol{\omega}; \mathbf{y})$ the mutual information between $\boldsymbol{\omega}$ and \mathbf{y} without explicit reference to the distribution.

Since \mathbf{x} is a sufficient statistic for \mathbf{y} , \mathbf{x} is more informative⁴ about \mathbf{y} than $\boldsymbol{\omega}$, that is $I(\boldsymbol{\omega}; \mathbf{y}) \leq I(\mathbf{x}; \mathbf{y})$. Then, the mutual information between the input and the output is no more than $C(Q)$ from the definition of the channel capacity.

The proof for general n and k is an elaboration of this argument. Since states are i.i.d., we can prove that for any strategy, the average payoff is the one induced by some splitting whose mutual information is no more than $\frac{k}{n} C(Q)$. The trick is to introduce an auxiliary random

⁴See [Cover and Thomas, 2006](#), Theorem 2.8.1, p. 34.

variable \mathbf{t} uniformly distributed over $\{1, \dots, n\}$ and known by the receiver. Then, we regard the average payoff over stages $1, \dots, n$ as the expected payoff for the randomly selected stage.

Second point, asymptotic construction. To make the intuition simple, let us consider a sequence of pairs of integers $(k_j, n_j)_{j \in \mathbb{N}}$ such that $k_j = n_j$ and let $k = n$ be a large term of this sequence. Take a splitting $(\lambda_m, \nu_m)_m$ of the prior μ which satisfies the information constraint. We want to show that for large n , there is a strategy σ of the sender such that for any best reply $\tau \in BR(\sigma)$ of the receiver, the payoff of the sender is at least about $\sum_m \lambda_m u_S^*(\nu_m)$. Let also $a_m^* \in A^*(\nu_m)$ such that $u_S^*(\nu_m) = \sum_{\omega} \nu_m(\omega) u_S(\omega, a_m^*)$.

A first intuition for the construction is as follows. From Shannon's coding Theorem⁵, if $I(\omega; \mathbf{m}) < C(Q)$, then for large n , there exists functions $f_1 : \Omega^n \rightarrow M^n$, $f_2 : M^n \rightarrow X^n$ and $g : Y^n \rightarrow M^n$, altogether a coding/decoding scheme, with the following properties. Given a sequence of states ω^n , the sender calculates a sequence of messages $m^n = f_1(\omega^n)$ such that with probability close to one, the empirical frequency of the (ω_t, m_t) 's is approximately the theoretical one $\lambda_m \nu_m(\omega)$. The sender then calculates a sequence of inputs $x^n = f_2(m^n)$ and sends them into the channel. If the receiver calculates $\hat{m}^n = g(y^n)$, then the messages are recovered with probability close to one: $\mathbb{P}(m^n = \hat{m}^n) \approx 1$.

This argument is standard in information theory but is not sufficient for proving our result. The proof is actually more complicated because the strategic receiver actually calculates the Bayesian posterior $\mathbb{P}(\omega^n | y^n)$ and chooses at stage t an action $a_t \in A^*(\mathbb{P}(\omega_t | y^n))$. Thus, the main task is to refine the construction in such a way that for any best reply of the receiver, with probability close to one, the optimal action $a_t \in A^*(\mathbb{P}(\omega_t | y^n))$ is equal to the recommended action $a_{\hat{m}_t}^*$ at most stages, that is, for a set of stages whose proportion is close to one. This implies that the payoff is approximately the target one.

The proof consists of three main steps. In the first step, we show that for each $\varepsilon > 0$, we can find a splitting ε -optimal for $V(\mu, c)$, which satisfies the information constraint with strict inequality and such that for each posterior ν_m , the action a_m^* which minimizes the sender payoff over $A^*(\nu_m)$ is *unique* in a neighborhood of ν_m . This latter property ensures that the receiver plays a_m^* whenever its belief is close to ν_m . We deduce that the difference between the realized payoff and the target payoff is bounded by the number of times t where the Bayesian posterior $\mathbb{P}(\omega_t | y^n)$ is far away from $\nu_{\hat{m}_t}$. The goal is then to show that this number is small with probability

⁵See Cover and Thomas, 2006, Theorem 10.4.1, p. 318.

close to one.

The second step consists in defining Shannon's strategy for this splitting. There, we adapt known construction from information theory to our setting.

At the third step, we prove that, under our construction, with probability close to one, the Bayesian posteriors $\mathbb{P}(\omega_t|y^n)$ are close enough to the target posteriors ν_{m_t} at most stages. This allows us to conclude that with probability close to one, the receiver plays the recommended actions at most stages and that the expected payoff is close to the target one. This step, where we estimate the realized Bayesian beliefs, is new compared to the information theoretic literature, which typically focuses on the average number of mistakes in decoding. Summing up, our construction is similar to the ones found in this literature but is adapted to the context where the receiver is maximizing its payoff. \square

3.1.2 Implications

We now provide some direct implications of the theorem.

Large capacity. Reordering the information constraint as $\sum_m \lambda_m H(\nu_m) \geq H(\mu) - c$, we see that if $c \geq H(\mu)$, the constraint is satisfied by all splittings. The value of the problem is thus the unconstrained concavification of u_S^* :

$$c \geq H(\mu) \Rightarrow V(\mu, c) = \text{cav } u_S^*(\mu).$$

As a consequence, if we fix n and Q and choose k large enough such that $\frac{k}{n}C(Q) \geq H(\mu)$, then the sender can achieve approximately the unconstrained maximum $\text{cav } u_S^*(\mu)$.

The intuition is simple: for fixed size of the state space, if the imperfect channel can be used a large number of times, then the sender is able to convey any message with arbitrarily high probability. More precisely, suppose $C(Q) > 0$ that is to say, $Q(\cdot|x)$ is not constant with respect to x . There exist distributions of inputs $p_m \in \Delta(X)$ that statistically identify the message:

$$m \neq m' \Rightarrow \sum_x p_m(x)Q(\cdot|x) \neq \sum_x p_{m'}(x)Q(\cdot|x).$$

For each message m , the sender can draw an i.i.d. sequence of messages x_1, \dots, x_k from p_m and sends them through the channel. The posterior belief of the receiver conditional on y_1, \dots, y_k then converges to the truth (the Dirac mass on m). Thus, asymptotically, the distributions of

actions of the receiver will be close to the one under perfect communication.

Small capacity. When c is close to 0, the information constraint $H(\mu) - \sum_m \lambda_m H(\nu_m) \leq c$ implies that the splitting is almost nonrevealing since:⁶

$$\sum_m \lambda_m \|\nu_m - \mu\|_1 \leq \sqrt{2 \ln 2 (H(\mu) - \sum_m \lambda_m H(\nu_m))}.$$

It follows that $V(\mu, c)$ is approximately $u_S^*(\mu)$, the payoff obtained without any information transmission.

As a consequence, if we fix Q and k , then for large n , the sender cannot get substantially more than $u_S^*(\mu)$.

Perfect channels. Our result applies to communication channels without noise. A communication channel has two sources of imperfection: the noise and the number of available messages, which is given exogenously. One insight of our work is that all that matters for the analysis is the capacity of the channel.

A channel (X, Y, Q) is called *perfect* if $X = Y$ and $Q(y|x) = \mathbb{1}_{\{x=y\}}$. For each integer $m \geq 2$, we denote Q_m^* the perfect communication channel with m messages where $m = |X| = |Y|$. Its capacity is⁷ $C(Q_m^*) = \log m$. We apply our results to the optimal robust payoff $U_S^*(\mu^n, Q_m^*)$ of the game where the persuasion problem is repeated n times and where the sender can send *one* message from a set with cardinality m . Our method applies since for large m , the channel Q_m^* can be seen as having the use of a binary perfect channel k times, with $k = \log_2 m$.

There are two simple extreme cases. First, if $m = 1$, the capacity of the channel is 0 and the sender cannot convey any information. Thus, $U_S^*(\mu^n, Q_m^*) = V(\mu, 0) = u_S^*(\mu)$. Second, if $m \geq |\Omega|^n$, then the sender can secure the unconstrained persuasion payoff $U_S^*(\mu^n, Q_m^*) = V(\mu, \log |\Omega|) = \text{cav } u_S^*(\mu)$ by treating each of the n problems separately and getting the payoff $\text{cav } u_S^*(\mu)$ for each instance. The first point of Theorem 3.1 shows that this is the best possible payoff.

More generally, Theorem 3.1 implies the following.

Corollary 3.2. *Consider a persuasion problem repeated n times, where the sender sends one message from a set of cardinality m . Then:*

⁶See Cover and Thomas, 2006, Lemma 11.6.1, p. 370.

⁷See Cover and Thomas, 2006, p. 184.

1. $U_S^*(\mu^n, Q_m^*) \leq V(\mu, \frac{\log m}{n})$.

2. For any pair of sequences of integers $(m_j, n_j)_{j \in \mathbb{N}}$ such that $\lim_{j \rightarrow \infty} \max(m_j, n_j) = \infty$ and $\lim_{j \rightarrow \infty} \frac{\log m_j}{n_j} = c$, we have $\lim_{j \rightarrow \infty} U_S^*(\mu^{n_j}, Q_{m_j}^*) = V(\mu, c)$.

Proof. The first point follows directly from Theorem 3.1. To see the second point, it is enough to remark that a perfect channel Q_m^* is “close” to k copies of a perfect binary channel with k such that $2^k \leq m < 2^{k+1}$, that is $k = \lfloor \log m \rfloor$. Having more messages at disposal is beneficial for the sender and thus $U_S^*(\mu^n, Q_m^*)$ is weakly increasing with m . It follows that:

$$U_S^*(\mu^n, (Q_2^*)^k) \leq U_S^*(\mu^n, Q_m^*) \leq U_S^*(\mu^n, (Q_2^*)^{k+1}).$$

Take a sequence $(m_j, n_j)_{j \in \mathbb{N}}$ such that $\lim_{j \rightarrow \infty} \max(m_j, n_j) = \infty$ and $\lim_{j \rightarrow \infty} \frac{\log m_j}{n_j} = c$, and define $k_j = \lfloor \log m_j \rfloor$. We have $\lim_{j \rightarrow \infty} \max(k_j, n_j) = \infty$, $\lim_{j \rightarrow \infty} \frac{k_j}{n_j} = c$ and the conclusion follows from Theorem 3.1.

3.2 Concavification with information constraint

In this section, we give some properties of the optimal splitting problem under information constraint. The motivation for this part of the results is two-fold. First, it is known that in a concavification problem, the number of posteriors (or of messages) can be chosen less than or equal to the number of states. One might wonder whether this remains true when there is a constraint on the feasible splittings. Second, models with costly information often use the mutual information as information cost (See e.g. Sims, 2003). We will see that in our case, this is derived by writing a Lagrangian for $V(\mu, c)$.

Consider the optimal splitting under information constraint:

$$\sup \left\{ \sum_m \lambda_m u_S^*(\nu_m) : \sum_m \lambda_m \nu_m = \mu, \sum_m \lambda_m H(\nu_m) \geq H(\mu) - c \right\}.$$

This is a special instance of the following optimization problem. Let $f, g : X \rightarrow \mathbb{R} \cup \{-\infty\}$ be two functions defined on a convex set $X \subseteq \mathbb{R}^d$, where X represents an abstract set of posteriors, f is a payoff function and g is a constraint capturing the feasible splittings. For $x \in X$ and $\gamma \in \mathbb{R}$ consider the problem:

$$F^g(x, \gamma) := \sup \left\{ \sum_m \lambda_m f(x_m) : \sum_m \lambda_m x_m = x, \sum_m \lambda_m g(x_m) \geq \gamma \right\}.$$

Let $f^g : X \times \mathbb{R} \rightarrow \mathbb{R} \cup \{-\infty\}$ defined by:

$$f^g(x, \gamma) = \begin{cases} f(x) & \text{if } \gamma \leq g(x), \\ -\infty & \text{otherwise.} \end{cases}$$

Theorem 3.3. *Then, for each $(x, \gamma) \in X \times \mathbb{R}$,*

1. $F^g(x, \gamma) = \text{cav } f^g(x, \gamma)$.
2. $F^g(x, \gamma) = \inf_{t \geq 0} \left\{ \text{cav } (f + tg)(x) - t\gamma \right\}$.

Applying this result to the optimal splitting under information constraint, we get:

Corollary 3.4. *For each $\mu \in \Delta(\Omega)$ and $c \geq 0$,*

1. $V(\mu, c)$ is the concavification of the function $u_S^H : \Delta(\Omega) \times \mathbb{R} \rightarrow \mathbb{R}$ defined as:

$$u_S^H(\nu, \eta) = \begin{cases} u_S^*(\nu) & \text{if } \eta \leq H(\nu), \\ -\infty & \text{otherwise,} \end{cases}$$

calculated at $(\nu, \eta) = (\mu, H(\mu) - c)$.

2. $V(\mu, c) = \inf_{t \geq 0} \left\{ \text{cav } (u_S^* + tH)(\mu) - t(H(\mu) - c) \right\}$.

Since it might be useful in other contexts, Theorem 3.3 is stated for general functions rather than specifically for the entropy function. This result has recently been generalized by [Doval and Skreta \(2018\)](#) to splitting problems with several constraints. The first point of the theorem states that the concavification with constraint, is the concavification of a bivariate function where an additional variable is added for the constraint (many variables when there are many constraints, see [Doval and Skreta, 2018](#)). The second point states that a Lagrangian function can be introduced and that the concavification under constraint is the concavification of the Lagrangian for some multiplier. The proof is in the Appendix (A.2).

A direct implication of the second point of Corollary 3.4 is that there exists $t^* = t^*(\mu, c)$ such that:

$$V(\mu, c) = \text{cav } (u_S^* + t^*H)(\mu) - t^*(H(\mu) - c).$$

To see the existence of t^* , notice that $\text{cav}(u_S^* + tH)(\mu) - t(H(\mu) - c) \geq (u_S^* + tH)(\mu) - t(H(\mu) - c) = u_S^*(\mu) + tc$, which tends to $+\infty$ as $t \rightarrow +\infty$. Therefore, $t \mapsto \text{cav}(u_S^* + tH)(\mu) - t(H(\mu) - c)$ reaches a minimum at some t^* .

If $(\lambda_m^*, \nu_m^*)_m$ is an optimal splitting, let $\mathcal{I}^* = H(\mu) - \sum_m \lambda_m^* H(\nu_m^*)$ be its mutual information. We have the following:

$$V(\mu, c) = \sum_m \lambda_m^* u_S^*(\nu_m^*) - t^*(\mathcal{I}^* - c). \quad (2)$$

We then find the usual Kuhn-Tucker slackness conditions. If $\mathcal{I}^* < c$, then $t^* = 0$ and the unconstrained optimum is feasible. If $t^* > 0$, the constraint is binding. The Lagrange multiplier t^* can be interpreted as the *shadow price of capacity*, that is, the marginal value of an extra unit of communication capacity.

This characterization can be related with the *cost of information* considered in the literature on rational inattention (See [Sims, 2003](#)) where the agent pays a cost proportional to the mutual information between the state and the signal he observes. In particular, [Caplin and Dean \(2013\)](#) consider the concavification of a utility function net of such an information cost. For persuasion games, [Gentzkow and Kamenica \(2014\)](#) assume that the sender pays a cost for choosing a disclosure strategy which is also related to the mutual information and also take the concavification of the net utility function.

Equation (2) can be seen as a microfoundation of the use of mutual information as the information cost: the limit optimal value of persuasion for a large number of copies of problems with communication over an imperfect channel, has the same value as a problem of persuasion with an information cost. There are some differences, however. First, the information cost is not the mutual information, but the difference between the mutual information and the capacity of the channel. That is, a cost reduces the payoff only when the sender would like to send more information bits than the capacity. Second, the unit price of capacity is endogenous and given by the Lagrange multiplier of the information constraint.

A direct implication is an upper bound of the number of posteriors needed to achieve the concavification.

Corollary 3.5. *In the optimization problem,*

$$V(\mu, c) = \sup \left\{ \sum_m \lambda_m u_S^*(\nu_m) : \sum_m \lambda_m \nu_m = \mu, \sum_m \lambda_m H(\nu_m) \geq H(\mu) - c \right\},$$

the number of posteriors can be chosen to be at most $\min\{|A|, |\Omega| + 1\}$.

Without the information constraint, the usual bound is $\min\{|A|, |\Omega|\}$: the number of posteriors or of messages can be upper bounded by the number of actions and the number of states. For the number of actions, the argument is that two messages for which the receiver chooses the same action can be merged into one and the corresponding two posteriors replaced by the average. The argument still holds due to the concavity of the entropy function: replacing two posteriors by their average increases the expected entropy and thus helps in satisfying the information constraint.

For the bound given by the number of states, the usual technical argument is that any point in the convex hull of the hypograph of a function on $\Delta(\Omega)$ is a convex combination involving $|\Omega|$ points. From Corollary 3.4, we consider the concavification of a function defined on $\Delta(\Omega) \times \mathbb{R}$ a domain with one extra dimension; thus, an extra posterior might be needed. A similar observation is made in Boleslavsky and Kim (2018), where due to an incentive constraint, an extra posterior is needed. In Section 4, we provide an example where $|\Omega| + 1$ posteriors are used at the optimum.

4 Illustrating example

4.1 Unrestricted communication

In this example, the sender is a firm that persuades the receiver to invest in a risky project. If the receiver does not invest (action a_0), the payoff is 0 for both players. If the receiver invests (action a_1), the project has return -7 in the bad state ω_0 and $+1$ in the good state ω_1 . Both states are equally likely. The sender receives a fee of $+1$ only if the receiver invests. The payoff table is as follows, the entries are pairs of payoffs for the players $i = S, R$ depending on the state and action.

	a_0	a_1	μ
ω_0	0, 0	1, -7	$\frac{1}{2}$
ω_1	0, 0	1, 1	$\frac{1}{2}$

The receiver invests for sure only when he holds a belief ν such that $\nu(\omega_1) > 7/8$. If $\nu(\omega_1) = 7/8$ he is indifferent. Assuming that in case of indifference he does not invest, the robust payoff of the sender $u_S^*(\nu)$ is 1 if $\nu(\omega_1) > 7/8$ and 0 otherwise.

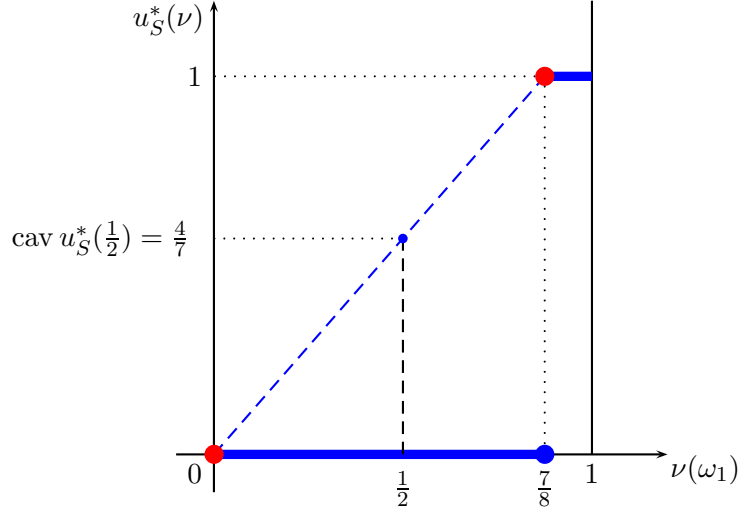


Figure 1: Concavification.

The concavification function $\text{cav } u_S^*(\nu)$ is continuous and equal to $\frac{8}{7}\nu(\omega_1)$ for $\nu(\omega_1) \leq \frac{7}{8}$ and 1 otherwise. It is easy to see that it does not depend on the action chosen by the receiver at $\nu(\omega_1) = \frac{7}{8}$, see Figure 1. If the receiver were to choose a_1 at the point of indifference, then the optimal splitting for the sender would be as follows:

$$\left(\frac{1}{2}, \frac{1}{2}\right) = \frac{3}{7}(1, 0) + \frac{4}{7}\left(\frac{1}{8}, \frac{7}{8}\right),$$

where a belief is denoted $\nu = (\nu(\omega_0), \nu(\omega_1))$. This yields a payoff of $\frac{4}{7}$ which is the highest that the sender can achieve given the uniform prior. For any small $\varepsilon > 0$, we can perturb the previous splitting and get the following:

$$\left(\frac{1}{2}, \frac{1}{2}\right) = \frac{3+8\varepsilon}{7+8\varepsilon}(1, 0) + \frac{4}{7+8\varepsilon}\left(\frac{1}{8} - \varepsilon, \frac{7}{8} + \varepsilon\right),$$

which achieves the payoff $\frac{4}{7+8\varepsilon}$ irrespective of the tie-breaking rule. Letting ε tend 0, we see that the sender achieves a payoff arbitrarily close to $\frac{4}{7}$, which is the optimal robust payoff.

4.2 Restricted and noisy communication

We consider binary sets of messages $X = \{x_0, x_1\}$, $Y = \{y_0, y_1\}$ and we assume that the channel has a *noise level* $\varepsilon \in [0, \frac{1}{2}]$, that is $Q(y_j|x_i) = \varepsilon$ for $j \neq i$, see Figure 2. The generic case is $\varepsilon \in (0, \frac{1}{2})$ where the label of the message (0 or 1) is changed with positive probability but observing a label 1 is still more likely when the input label is 1. When $\varepsilon = \frac{1}{2}$, the distribution of the output message is independent from the input message, so the channel completely disrupts the communication.

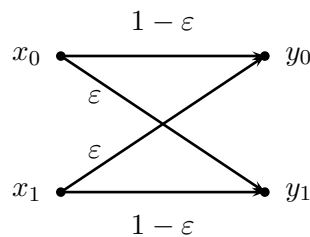


Figure 2: Binary symmetric channel.

A special case is the *binary perfect channel* when $\varepsilon = 0$: identifying together the sets X and Y , an input message x is received with certainty. Communication is then restricted only by the number of available messages, i.e. the cardinality of X .

The capacity of the binary symmetric channel⁸ is $1 - H(\varepsilon)$ where with some abuse of notation, $H(\varepsilon)$ denotes the entropy of the binary probability distribution $(\varepsilon, 1 - \varepsilon)$.

4.3 One-shot scenario $k = n = 1$

Let a strategy σ of the sender be parametrized by $\sigma(x_0|\omega_0) = 1 - \alpha$ and $\sigma(x_1|\omega_1) = 1 - \beta$; see Figure 3.

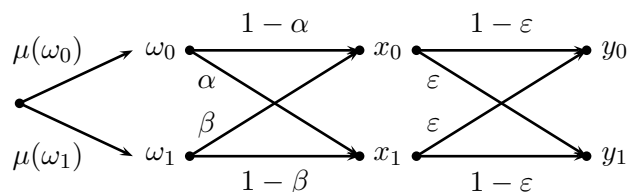


Figure 3: Strategy on the binary symmetric channel.

⁸Cover and Thomas, 2006, Example 2.1.1, p. 15.

Then, $\mathbb{P}_\sigma(y_1|\omega_0) = \alpha(1 - \varepsilon) + (1 - \alpha)\varepsilon$, $\mathbb{P}_\sigma(y_0|\omega_1) = \beta(1 - \varepsilon) + (1 - \beta)\varepsilon$ and from Bayes' rule,

$$\mathbb{P}_\sigma(\omega_1|y_1) = \frac{\mu(\omega_1)(1 - \mathbb{P}_\sigma(y_0|\omega_1))}{\mu(\omega_0)\mathbb{P}_\sigma(y_1|\omega_0) + \mu(\omega_1)(1 - \mathbb{P}_\sigma(y_0|\omega_1))},$$

$$\mathbb{P}_\sigma(\omega_1|y_0) = \frac{\mu(\omega_1)\mathbb{P}_\sigma(y_0|\omega_1)}{\mu(\omega_0)(1 - \mathbb{P}_\sigma(y_1|\omega_0)) + \mu(\omega_1)\mathbb{P}_\sigma(y_0|\omega_1)}.$$

It is easy to see that the numbers $\mathbb{P}_\sigma(y_1|\omega_0)$, $\mathbb{P}_\sigma(y_0|\omega_1)$, $\mathbb{P}_\sigma(\omega_1|y_1)$, $\mathbb{P}_\sigma(\omega_1|y_0)$ all belong to the interval $[\varepsilon, 1 - \varepsilon]$.

A pair of posteriors (ν_0, ν_1) is said to be *feasible* in the one-shot scenario if there exists a number $\lambda \in [0, 1]$ such that:

$$(\mu(\omega_0), \mu(\omega_1)) = \lambda(\nu_0(\omega_0), \nu_0(\omega_1)) + (1 - \lambda)(\nu_1(\omega_0), \nu_1(\omega_1)).$$

The feasible splittings can be characterized as follows.

Lemma 4.1. *We consider the one-shot problem where $n = k = 1$. A pair of posteriors (ν_0, ν_1) is feasible if and only if $\nu_1 = \nu_0 = \mu$ or,*

$$\varepsilon \leq \frac{\nu_0(\omega_1)(\nu_1(\omega_1) - \mu(\omega_1))}{\mu(\omega_1)(\nu_1(\omega_1) - \nu_0(\omega_1))} \leq 1 - \varepsilon$$

and

$$\varepsilon \leq \frac{(1 - \nu_0(\omega_1))(\mu(\omega_1) - \nu_0(\omega_1))}{(1 - \mu(\omega_1))(\nu_1(\omega_1) - \nu_0(\omega_1))} \leq 1 - \varepsilon.$$

The proof is in Appendix A.1. As an illustration, take the uniform prior $(\frac{1}{2}, \frac{1}{2})$ and a level of noise $\varepsilon = \frac{1}{4}$. The feasible posteriors are shown by the colored green regions on Figure 5.

From the previous discussion, it is impossible to induce beliefs with $\nu(\omega_1) > \frac{3}{4}$. Therefore, the receiver will never be confident enough to invest and the payoff is 0 for the sender.

4.4 Asymptotic scenario with $k = n \rightarrow \infty$

We consider the case where $k = n$ tends to infinity with a noise level of $\varepsilon = \frac{1}{4}$ and compute the value of the optimal splitting problem with information constraint. The capacity of the channel is $1 - H(\frac{1}{4})$, the entropy of the uniform prior is 1; therefore, the information constraint is $\sum_m \lambda_m H(\mu_m) \geq H(\frac{1}{4})$. Under this constraint the optimal splitting for the sender satisfies:

$$\left(\frac{1}{2}, \frac{1}{2}\right) = \lambda \left(\frac{1}{8}, \frac{7}{8}\right) + (1 - \lambda)(\nu_0(\omega_0), \nu_0(\omega_1))$$

and

$$H\left(\frac{1}{4}\right) = \lambda H\left(\frac{7}{8}\right) + (1 - \lambda)H(\nu_0(\omega_1)).$$

To see why it is optimal, first observe that the sender has to bring on some posterior, denoted by ν_1 , with $\nu_1(\omega_1) > \frac{7}{8}$ in order to get some payoff. To get it with the highest probability, he should aim for $\nu_1(\omega_1) = \frac{7}{8}$. Among the posteriors that induce investment, this is also the one with highest entropy. Second, to maximize expected payoffs, the remaining posteriors must be as far away as possible from the prior; that is, the information constraint should bind. Additionally, note that only one posterior, denoted by ν_0 , will be optimally generated in the region $\nu_0(\omega_1) < \frac{7}{8}$. Since the entropy is strictly concave, replacing two posteriors on this region by their average does not change the payoff and increases the entropy.

Solving these two equations numerically we get, $\nu_0(\omega_1) \approx 0.340$ and $V(\mu, Q) = \lambda \approx 0.298$ instead of the zero value for the one-shot scenario and about 52.1% of the unconstrained optimum $\frac{4}{7}$.

This is shown in Figure 4 which plots the payoff function and the entropy function. The splitting of μ into ν_0, ν_1 is shown by the three points on the horizontal axis. On the vertical line $\mu(\omega_1) = \frac{1}{2}$, we can read the average payoff with the red line and the average entropy with the green line. To see optimality on the picture, if we move $\nu_0(\omega_1)$ to the right, then the average payoff decrease, and if we move it to the left, the average entropy will fall below $H(\frac{1}{4})$ and the information constraint will be violated.

The optimal splitting is also marked on Figure 5 which shows the set of pairs of posteriors for the splittings that satisfy the information constraint (union of green and blue regions).

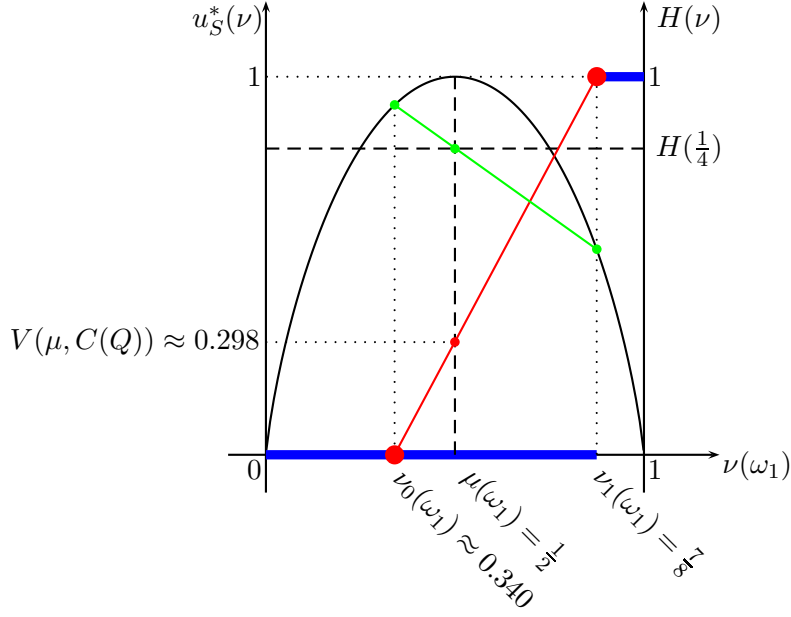


Figure 4: For a noise parameter $\varepsilon = \frac{1}{4}$, the optimal splitting is given by $\nu_0(\omega_1) \approx 0.340$ and $\nu_1(\omega_1) = \frac{7}{8}$.

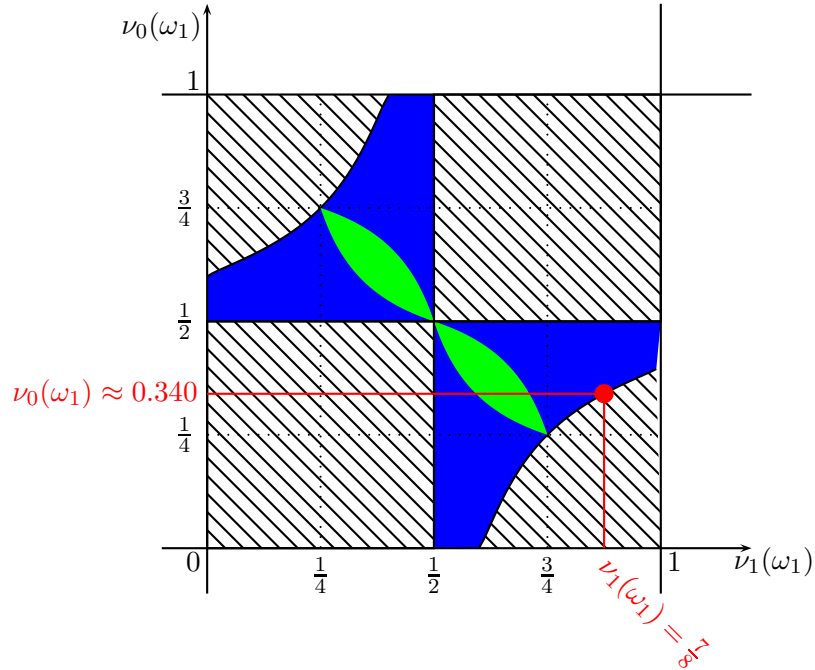


Figure 5: For a noise parameter $\varepsilon = \frac{1}{4}$, the green lenses correspond to the feasible posteriors (ν_0, ν_1) characterized in Lemma 4.1 for the one-shot scenario $k = n = 1$. The blue and green regions correspond to the feasible posteriors (ν_0, ν_1) in the asymptotic scenario where $k = n \rightarrow \infty$. The red point corresponds to the optimal splitting, also depicted in Figure 4. The hatched areas correspond to the nonfeasible posteriors (ν_0, ν_1) .

On Figure 6, we represent the value $V(\mu, C(Q))$ of the optimal splitting problem as a function of the prior μ , for different values for the noise parameter $\varepsilon \in \left\{ \frac{1}{20}, \frac{3}{20}, \frac{1}{4}, \frac{7}{20}, \frac{9}{20}, \frac{99}{200} \right\}$. It is found by solving the following system for ν_0 :

$$(\mu(\omega_0), \mu(\omega_1)) = \lambda \left(\frac{1}{8}, \frac{7}{8} \right) + (1 - \lambda)(\nu_0(\omega_0), \nu_0(\omega_1))$$

and

$$H(\mu(\omega_1)) - 1 + H(\varepsilon) = \lambda H\left(\frac{7}{8}\right) + (1 - \lambda)H(\nu_0(\omega_1)).$$

When $\mu(\omega_1) = \frac{1}{2}$ and $\varepsilon = \frac{1}{4}$, we recover the value $V(\mu, C(Q)) \approx 0.298$ as in Figure 4.

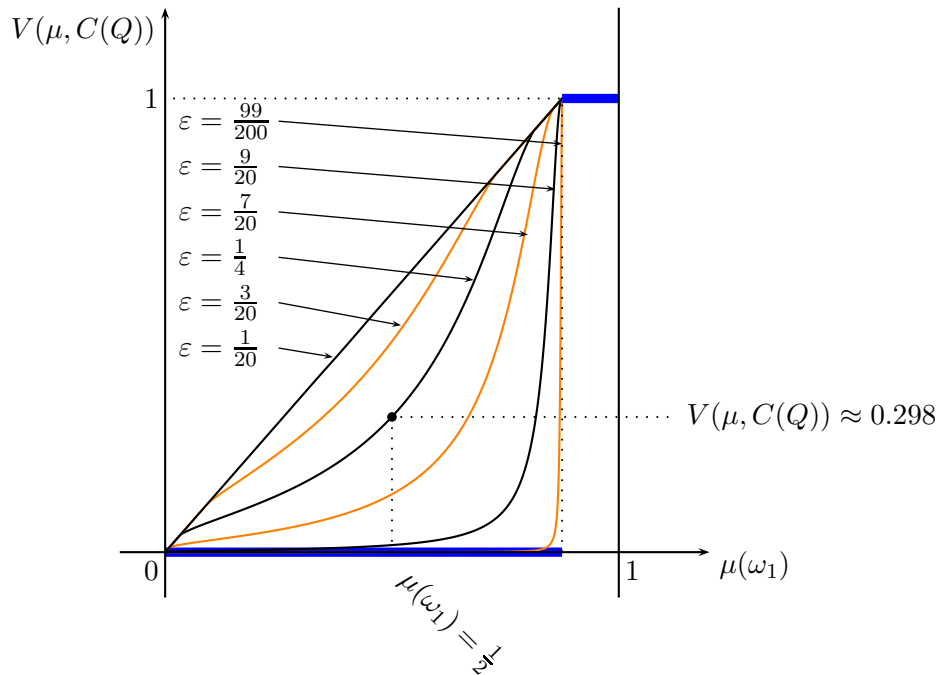


Figure 6: Value of the optimal splitting problem as a function of the prior μ , for different noise parameters $\varepsilon \in \left\{ \frac{1}{20}, \frac{3}{20}, \frac{1}{4}, \frac{7}{20}, \frac{9}{20}, \frac{99}{200} \right\}$.

Observe that the function $V(\mu, C(Q))$ is *not* concave with respect to the prior μ . From Corollary 3.4, $V(\mu, c)$ is the concavification of the function u_S^H calculated at $(\mu, H(\mu) - c)$, so this composed function need not be concave.

4.5 Perfect binary channel with $k = \frac{n}{2} \rightarrow \infty$

We consider the same example as before, repeated n times with the uniform prior $\mu = (\frac{1}{2}, \frac{1}{2})$. In line with the motivating example from the introduction, we consider a perfect channel and assume that the sender has at its disposal half as many messages as needed to communicate

perfectly, that is $k = \frac{n}{2}$. Since the capacity of the binary perfect channel is one, $\frac{k}{n} = \frac{1}{2}$ and $H(\mu) = 1$, the information constraint is:

$$H(\mu) - \sum_m \lambda_m H(\mu_m) \leq \frac{1}{2} \iff \sum_m \lambda_m H(\mu_m) \geq \frac{1}{2}.$$

Observe that this constraint is identical to the one obtained with a binary symmetric channel with noise ε such that $H(\varepsilon) = \frac{1}{2}$ (i.e., $\varepsilon \approx 0.110$). Therefore, the optimal splitting is given by the following system:

$$(\mu(\omega_0), \mu(\omega_1)) = \lambda \left(\frac{1}{8}, \frac{7}{8} \right) + (1 - \lambda)(\nu_0(\omega_0), \nu_0(\omega_1))$$

and

$$\frac{1}{2} = \lambda H\left(\frac{7}{8}\right) + (1 - \lambda)H(\nu_0(\omega_1)).$$

Solving numerically, we find $V(\mu, \frac{1}{2}) \approx 0.519$.

4.6 On the number of posteriors

We give now an example showing the tightness of the bound $\min\{|A|, |\Omega| + 1\}$ on the number of posteriors, given in Corollary 3.5. The payoff table is as follows:

	a_0	a_1	a_2	
ω_0	0, 0	1, -7	1, 1	$\frac{1}{2}$
ω_1	0, 0	1, 1	1, -7	$\frac{1}{2}$

There are two risky projects (a_1 and a_2) and the sender wants to persuade the receiver to invest in any of them. The receiver invests only if $\nu(\omega_1) > 7/8$ or $\nu(\omega_1) < 1/8$.

With unrestricted communication, the solution is clear: the sender fully discloses the state and gets a payoff of 1. However, with a binary symmetric channel with noise $\varepsilon = 1/4$, the sender gets 0 in the one-shot scenario. Consider now the case where $k = n \rightarrow \infty$.

The “one-sided” solution of Section 4.4 is feasible. Recall that this is the splitting such that:

$$\left(\frac{1}{2}, \frac{1}{2} \right) = \lambda \left(\frac{1}{8}, \frac{7}{8} \right) + (1 - \lambda)(\nu_0(\omega_0), \nu_0(\omega_1))$$

and

$$H\left(\frac{1}{4}\right) = \lambda H\left(\frac{7}{8}\right) + (1 - \lambda)H(\nu_0(\omega_0), \nu_0(\omega_1)).$$

with $\nu_0(\omega_1) \approx 0.340$ and $\lambda \approx 0.298$. It is easy to see that this is optimal among the splittings with two posteriors. Indeed, it is not possible that the two posteriors induce investment while satisfying the information constraint.

However, this is not optimal. The optimal splitting has three posteriors and is as follows:

$$\left(\frac{1}{2}, \frac{1}{2}\right) = (1 - \lambda)\left(\frac{1}{2}, \frac{1}{2}\right) + \frac{\lambda}{2}\left(\frac{1}{8}, \frac{7}{8}\right) + \frac{\lambda}{2}\left(\frac{7}{8}, \frac{1}{8}\right)$$

with

$$H\left(\frac{1}{4}\right) = (1 - \lambda)H\left(\frac{1}{2}\right) + \frac{\lambda}{2}H\left(\frac{1}{8}\right) + \frac{\lambda}{2}H\left(\frac{7}{8}\right).$$

This pins down a unique λ and solving numerically yields $\lambda \approx 0.413$. Since λ is the probability of investment, we get $V(\mu, Q) \approx 0.413$ which is about 38% better than what is achieved with a splitting with two points.

To see that this is optimal, first since there are two states, we know that three posteriors are sufficient. Second, it is not possible to have all posteriors in the investment region and to satisfy the information constraint. If there is only one posterior in the investment region, then the splitting achieves no more than the “one-sided” solution. Therefore, it is optimal to have two posteriors in the investment region and one outside of it. However, then, it is optimal to choose the point in the middle region to be $(\frac{1}{2}, \frac{1}{2})$, since this is the one with the highest entropy.

5 Beyond identical problems

The main result can be extended to series of persuasion problems which are not all identical, but such that each type of problem is repeated many times. Suppose that we have a family of persuasion problems indexed by a type parameter z in a finite set Z . That is, for every $z \in Z$, there is a prior probability distribution $\mu(\cdot|z) \in \Delta(\Omega)$ and payoff functions $u_i(\cdot, z) : \Omega \times A \rightarrow \mathbb{R}$ for each player $i = S, R$. The series of persuasion problems is given by a sequence $z^n = (z_1, \dots, z_n)$ which is commonly known by both players. The distribution of states is as follows:

$$\mu^n(\omega^n|z^n) := \prod_{t=1}^n \mu(\omega_t|z_t).$$

If the sequence of states and actions are respectively ω^n, a^n , the payoff for player i is $\frac{1}{n} \sum_{t=1}^n u_i(\omega_t, a_t, z_t)$.

The communication technology is still given by a channel $Q : X \rightarrow \Delta(Y)$ used k times, so that the strategy sets are the same as before for both players. The optimal robust payoff of the sender is defined as before and is denoted by $U_S^*(\mu^n, Q^k, z^n)$.

For each posterior belief $\nu \in \Delta(\Omega)$ and type $z \in Z$, the set of optimal actions of the receiver is $A^{*z}(\nu) = \operatorname{argmax}_{a \in A} \sum_{\omega} \nu(\omega) u_R(\omega, a, z)$ and we denote by $u_S^{*z}(\nu) = \min_{a \in A^{*z}(\nu)} \sum_{\omega} \nu(\omega) u_S(\omega, a, z)$ the robust payoff of the sender at the belief ν .

Definition 5.1. For $\pi \in \Delta(Z)$ and $c \geq 0$, the optimal splitting problem with information constraint is as follows:

$$\begin{aligned} V^Z(\mu, c, \pi) = & \sup \sum_z \pi(z) \sum_m \lambda_m^z u_S^{*z}(\nu_m^z) \\ \text{s.t.} & \sum_m \lambda_m^z \nu_m^z = \mu(\cdot|z), \quad \forall z \in Z, \\ \text{and} & \sum_z \pi(z) \left(H(\mu(\cdot|z)) - \sum_m \lambda_m^z H(\nu_m^z) \right) \leq c. \end{aligned}$$

The interpretation is as follows. Suppose that $\pi(z)$ represents the probability, or frequency, of occurrence of z . Conditional on z which is known by both players, the sender performs a spitting of $\mu(\cdot|z)$, $\sum_m \lambda_m^z \nu_m^z = \mu(\cdot|z)$, and gets the payoff $\sum_m \lambda_m^z u_S^{*z}(\nu_m^z)$. The information constraint imposes the *average* mutual information to be less than or equal to the capacity.

Given a sequence $z^n \in Z^n$, let $\pi_n \in \Delta(Z)$ be the empirical frequency induced by the sequence: for each $z \in Z$, $\pi_n(z) = \frac{1}{n} |\{t : z_t = z\}|$.

Theorem 5.2. 1. The optimal robust payoff of the sender is no more than the value of the optimal splitting problem with information constraint. For each pair of integers n, k :

$$U_S^*(\mu^n, Q^k, z^n) \leq V^Z\left(\mu, \frac{k}{n} C(Q), \pi_n\right).$$

2. The optimal robust payoff of the sender converges to the value of the optimal splitting problem with information constraint in the following sense. For each $\pi \in \Delta(Z)$ and $r \in [0, +\infty]$, for each pair of sequences of integers $(k_j, n_j)_{j \in \mathbb{N}}$ such that $\lim_{j \rightarrow \infty} \max(n_j, k_j) = \infty$, $\lim_{j \rightarrow \infty} \frac{k_j}{n_j} = r$ and $\lim_{j \rightarrow \infty} \pi_{n_j} = \pi$, we have:

$$\lim_{j \rightarrow \infty} U_S^*(\mu^{n_j}, Q^{k_j}, z^{n_j}) = V^Z(\mu, rC(Q), \pi).$$

Given a sequence z^n , π_n is the empirical distribution of types of problems. The optimal payoff of the sender is bounded above by the value of optimal splitting under information constraint. Suppose that the distribution of types is held fixed (or converges to) π , then when n and k grow large, the sender is able to secure approximately this value. The arguments of the proof of Theorem 3.1 extend quite easily to this case (up to lengthy adaptations for the second point) so the proof is omitted.

This extension applies to the case where the proportions of types of problems are fixed. Alternatively, the sequence z^n could be drawn i.i.d. from a prior distribution $\pi \in \Delta(Z)$.

Notice the channel Q^k is used for transferring information about all problems. Thus, Theorem 5.2 does more than merely patching up distinct families of problems together. The capacity of the channel bounds the total amount of information, across all problems. Thus, all problems, even of different types, are linked together in the messages.

6 Conclusion

We have analyzed a persuasion game where the sender communicates with the receiver through a fixed and imperfect channel. The optimal payoff of the sender is bounded above by the value of the optimal splitting problem with information constraint. When the sender and the receiver are engaged in many repetitions of identical persuasion games, the optimal payoff for the sender converges to the upper bound as the number of repetitions increases.

There are several interesting variations or extensions of this model.

1. *Private information of the receiver.* In the model, it is assumed that the information about the state is fully controlled by the sender. To model private information of the receiver, consider the extension of the previous section, let nature draw pairs $(\omega_t, z_t)_{t=1, \dots, n}$ and assume that ω^n is the private information of the sender and z^n is the private information of the receiver. Our methods generalize to this case provided that we use a suitable generalization of the information constraint. A random message \mathbf{m} can be transmitted over the channel provided that its mutual information with the state, *conditional on the private information of the receiver* $I(\boldsymbol{\omega}; \mathbf{m} | \mathbf{z}) \leq C$ is less than or equal to the capacity, where

$$I(\boldsymbol{\omega}; \mathbf{m} | \mathbf{z}) := \sum_z \mathbb{P}(z) I(\boldsymbol{\omega}; \mathbf{m} | \mathbf{z} = z)$$

is the expectation over \mathbf{z} of the mutual information conditional on $\{\mathbf{z} = z\}$.

2. *Commitment of the receiver.* In the persuasion model, the sender first chooses its strategy and is committed to playing it. A natural twist is to let the receiver choose his strategy first and commit to it. This turns into a mechanism design problem where the receiver is a principal offering a contract to an informed agent (the sender), and where the agent communicates with the principal through an imperfect channel. Again, an information constraint holds, but the impact of incentives is different. Namely for each sequence of states $(\omega_1, \dots, \omega_n)$, the sender/agent should not have an incentive to behave as if it was another one $(\omega'_1, \dots, \omega'_n)$. The task is to prove that using the usual coding scheme is indeed an optimal strategy for the agent, or rather that any optimal strategy is not too different from the coding scheme. This variation is under close study.

3. *More general processes.* It would be interesting to generalize the results to a larger class of stochastic processes of states. Information theoretic methods can be extended to Markov chains, see [Cover and Thomas \(2006\)](#) and to more general processes, see [Han \(2003\)](#). While an information constraint would certainly hold, it is an open problem to characterize the optimal payoff for the sender. What is the best way to exploit the correlations between states?

References

- AKYOL, E., C. LANGBORT, AND T. BAŞAR (2017): “Information-Theoretic Approach to Strategic Communication as a Hierarchical Game,” *Proceedings of the IEEE*, 105(2), 205–218.
- AUMANN, R., AND M. MASCHLER (1995): *Repeated Games with Incomplete Information*. MIT Press, Cambridge, MA.
- BERGEMANN, D., AND S. MORRIS (2016): “Information Design, Bayesian Persuasion, and Bayes Correlated Equilibrium,” *American Economic Review Papers and Proceedings*, 106(5), 586–591.
- (2017): “Information Design: a Unified Perspective,” *Cowles Foundation Discussion Paper No 2075*.
- BLUME, A., O. J. BOARD, AND K. KAWAMURA (2007): “Noisy Talk,” *Theoretical Economics*, 2, 395–440.
- BOLES LAVSKY, R., AND C. COTTON (2015): “Grading Standards and Education Quality,” *American Economic Journal: Microeconomics*, 7(2), 248–279.

- BOLESZLAVSKY, R., AND K. KIM (2018): “Bayesian Persuasion and Moral Hazard,” *Working paper*.
- CAPLIN, A., AND M. DEAN (2013): “Behavioral Implications of Rational Inattention with Shannon Entropy,” *NBER Working Papers 19318*.
- COVER, T. M., AND J. A. THOMAS (2006): *Elements of Information Theory*. 2nd. Ed., Wiley-Interscience, New York.
- CRAWFORD, V. P., AND J. SOBEL (1982): “Strategic Information Transmission,” *Econometrica*, 50(6), 1431–1451.
- CUFF, P., H. PERMUTER, AND T. COVER (2010): “Coordination Capacity,” *IEEE Transactions on Information Theory*, 56(9), 4181–4206.
- CUFF, P., AND L. ZHAO (2011): “Coordination using Implicit Communication,” *Proceedings of the IEEE Information Theory Workshop (ITW)*, pp. 467–471.
- DOVAL, L., AND V. SKRETA (2018): “Constrained Information Design: Toolkit,” *Working paper*.
- GAMAL, A. E., AND Y.-H. KIM (2011): *Network Information Theory*. Cambridge University Press.
- GENTZKOW, M., AND E. KAMENICA (2014): “Costly Persuasion,” *American Economic Review*, 104, 457–462.
- GOSSNER, O., P. HERNÁNDEZ, AND A. NEYMAN (2006): “Optimal Use of Communication Resources,” *Econometrica*, 74(6), 1603–1636.
- GOSSNER, O., AND T. TOMALA (2006): “Empirical Distributions of Beliefs under Imperfect Observation,” *Mathematics of Operation Research*, 31(1), 13–30.
- (2007): “Secret Correlation in Repeated Games with Imperfect Monitoring,” *Mathematics of Operation Research*, 32(2), 413–424.
- GOSSNER, O., AND N. VIEILLE (2002): “How to Play with a Biased Coin?,” *Games and Economic Behavior*, 41(2), 206–226.
- HAN, T. S. (2003): *Information-spectrum Methods in Information Theory*. Springer.

- HEBERT, B., AND M. WOODFORD (2018): “Information Costs and Sequential Information Sampling,” *NBER Working Paper 25316*.
- HERNÁNDEZ, P., AND B. VON STENGEL (2014): “Nash Codes for Noisy Channels,” *Operations Research*, 62(6), 1221–1235.
- INOSTROZA, N., AND A. PAVAN (2018): “Persuasion in Global Games with Application to Stress Testing,” *working paper*.
- JACKSON, M. O., AND H. F. SONNENSCHNEIN (2007): “Overcoming Incentive Constraints by Linking Decisions,” *Econometrica*, 75(1), 241–257.
- KAMENICA, E., AND M. GENTZKOW (2011): “Bayesian Persuasion,” *American Economic Review*, 101, 2590–2615.
- LE TREUST, M. (2017): “Joint Empirical Coordination of Source and Channel,” *IEEE Transactions on Information Theory*, 63(8), 5087–5114.
- LE TREUST, M., AND T. TOMALA (2016): “Information Design for Strategic Coordination of Autonomous Devices with Non-Aligned Utilities,” *Proceedings of the IEEE 54th Allerton conference, Monticello, Illinois*, pp. 233–242.
- MARTIN, D. (2017): “Strategic Pricing with Rational Inattention to Quality,” *Games and Economic Behavior*, 104, 131–145.
- MATEJKA, F., AND A. MCKAY (2015): “Rational Inattention to Discrete Choices: A New Foundation for the Multinomial Logit Model,” *American Economic Review*, 105(1), 272–98.
- MATHEVET, L., J. PEREGO, AND I. TANEVA (2019): “On Information Design in Games,” *working paper, to appear in Journal of Political Economy*.
- MATYSKOVA, L. (2018): “Bayesian Persuasion With Costly Information Acquisition,” *Working Paper*.
- MERHAV, N., AND S. SHAMAI (2007): “Information Rates Subject to State Masking,” *IEEE Transactions on Information Theory*, 53(6), 2254–2261.
- MORRIS, S., AND P. STRACK (2019): “The Wald Problem and the Equivalence of Sequential Sampling and Static Information Costs,” *Working Paper*.

- NEYMAN, A., AND D. OKADA (1999): “Strategic Entropy and Complexity in Repeated Games,” *Games and Economic Behavior*, 29(1–2), 191–223.
- (2000): “Repeated Games with Bounded Entropy,” *Games and Economic Behavior*, 30(2), 228–247.
- PEREZ, E., AND V. SKRETA (2018): “Test Design under Falsification,” *working paper*.
- ROCKAFELLAR, R. (1970): *Convex Analysis*, Princeton landmarks in mathematics and physics. Princeton University Press.
- SHANNON, C. (1948): “A Mathematical Theory of Communication,” *Bell System Technical Journal*, 27, 379–423.
- (1959): “Coding Theorems for a Discrete Source with a Fidelity Criterion,” *IRE National Convention Record, Part 4*, pp. 142–163.
- SIMS, C. (2003): “Implication of Rational Inattention,” *Journal of Monetary Economics*, 50(3), 665–690.
- STEINER, J., C. STEWART, AND F. MATEJKA (2017): “Rational Inattention Dynamics: Inertia and Delay in Decision-Making,” *Econometrica*, 84(2), 521–553.
- TANEVA, I. (2018): “Information Design,” *working paper, to appear in American Economic Journal: Microeconomics*.
- TSAKAS, E., AND N. TSAKAS (2018): “Noisy Persuasion,” *Working Paper*.

A Appendix

This appendix contains all the formal proofs.

A.1 Proof of Lemma 4.1

For a, b in $[0, 1]$, consider the system:

$$\nu_1(\omega_1) = \frac{\mu(\omega_1)(1-b)}{\mu(\omega_0)a + \mu(\omega_1)(1-b)}, \quad \nu_0(\omega_1) = \frac{\mu(\omega_1)b}{\mu(\omega_0)(1-a) + \mu(\omega_1)b}. \quad (3)$$

If $\nu_1 = \nu_0 = \mu$, then it must be that $a = 1 - b$. Otherwise, $\nu_1(\omega_1) \neq \nu_0(\omega_1)$. It is easily verified that the system has a unique solution given by:

$$b = \frac{\nu_0(\omega_1)(\nu_1(\omega_1) - \mu(\omega_1))}{\mu(\omega_1)(\nu_1(\omega_1) - \nu_0(\omega_1))}$$

and

$$a = \frac{(1 - \nu_0(\omega_1))(\mu(\omega_1) - \nu_0(\omega_1))}{(1 - \mu(\omega_1))(\nu_1(\omega_1) - \nu_0(\omega_1))}.$$

Take a strategy σ defined by $\sigma(x_0|\omega_0) = 1 - \alpha$ and $\sigma(x_1|\omega_1) = 1 - \beta$ and a binary symmetric channel with noise ε . The posteriors ν_1, ν_0 are given by the system (3) for $a := \alpha(1 - \varepsilon) + \varepsilon(1 - \alpha)$ and $b := \beta(1 - \varepsilon) + \varepsilon(1 - \beta)$. As α, β vary in $[0, 1]$, a and b range freely over $[\varepsilon, 1 - \varepsilon]$,

$$\{(\alpha(1 - \varepsilon) + \varepsilon(1 - \alpha), \beta(1 - \varepsilon) + \varepsilon(1 - \beta)) : (\alpha, \beta) \in [0, 1]^2\} = [\varepsilon, 1 - \varepsilon]^2.$$

This concludes the proof. □

A.2 Proofs for Sections 3.2

A.2.1 Proof of Theorem 3.3, point 1

The function $\text{cav } f^g(x, \gamma)$ is given by the following program:

$$\begin{aligned} & \sup \quad \sum_m \lambda_m f(x_m) \\ & \text{s.t.} \quad \sum_m \lambda_m x_m = x, \quad \sum_m \lambda_m \gamma_m = \gamma \\ & \text{and} \quad \forall m, \gamma_m \leq g(x_m). \end{aligned}$$

Take a family $(\lambda_m, x_m, \gamma_m)_m$ feasible for this program. We have $\sum_m \lambda_m g(x_m) \geq \gamma$, thus this family is feasible for $F^g(x, \gamma)$. Therefore, $\text{cav } f^g(x, \gamma) \leq F^g(x, \gamma)$.

Conversely, take a family $(\lambda_m, x_m)_m$ such that $\sum_m \lambda_m x_m = x$ and $\sum_m \lambda_m g(x_m) \geq \gamma$. Let $\bar{\gamma} = \sum_m \lambda_m g(x_m)$ and for each m , $\gamma_m = g(x_m) + \gamma - \bar{\gamma}$. Then, $\sum_m \lambda_m \gamma_m = \gamma$ and since $\bar{\gamma} \geq \gamma$, for each m , $\gamma_m \leq g(x_m)$. Thus, $(\lambda_m, x_m, \gamma_m)_m$ is feasible for $\text{cav } f^g(x, \gamma)$ and $\text{cav } f^g(x, \gamma) \geq F^g(x, \gamma)$. □

A.2.2 Proof of Theorem 3.3, point 2

Recall that the Fenchel conjugate of $f : X \subseteq \mathbb{R}^d \rightarrow \mathbb{R}$ is $f^*(p) = \sup_x \{x \cdot p - f(x)\}$, where $x \cdot p$ denotes the inner product. Then, the largest convex function below f is equal to $(f^*)^*$ (Rockafellar, 1970, Corollary 12.1.1, p. 103), therefore $(f^*)^*(x) = -\text{cav}(-f)(x)$. Playing with signs, it follows that:

$$\text{cav } f(x) = \inf_p \left\{ x \cdot p + \sup_y \{f(y) - p \cdot y\} \right\}. \quad (4)$$

We apply this formula to the function:

$$f^g(x, \gamma) = \begin{cases} f(x) & \text{if } \gamma \leq g(x), \\ -\infty & \text{otherwise.} \end{cases}$$

This gives,

$$\begin{aligned} \text{cav } f^g(x, \gamma) &= \inf_{p, z} \left\{ p \cdot x + z\gamma + \sup_{y, \eta} \{f^g(y, \eta) - p \cdot y - z\eta\} \right\} \\ &= \inf_{p, z} \left\{ p \cdot x + z\gamma + \sup_{y, \eta: \eta \leq g(y)} \{f(y) - p \cdot y - z\eta\} \right\}. \end{aligned}$$

If $z > 0$ then by letting $\eta \rightarrow -\infty$, the sup is $+\infty$. Therefore, in the infimum we can restrict to $z \leq 0$. Setting $t = -z \geq 0$ we get:

$$\begin{aligned} \text{cav } f^g(x, \gamma) &= \inf_{t \geq 0, p} \left\{ p \cdot x - t\gamma + \sup_{y, \eta: \eta \leq g(y)} \{f(y) - p \cdot y + t\eta\} \right\} \\ &= \inf_{t \geq 0, p} \left\{ p \cdot x - t\gamma + \sup_y \{f(y) - p \cdot y + tg(y)\} \right\} \\ &= \inf_{t \geq 0} \left\{ \inf_p \left\{ p \cdot x + \sup_y \{f(y) + tg(y) - p \cdot y\} \right\} - t\gamma \right\} \end{aligned}$$

where the second line holds since $t \geq 0$ and the third line is just reorganizing. The result follows by remarking that $\inf_p \{p \cdot x + \sup_y \{f(y) + tg(y) - p \cdot y\}\} = \text{cav}(f + tg)(x)$. \square

A.2.3 Proof of Corollary 3.5, upper bound $|\Omega| + 1$

Corollary 3.5 follows from a well-known fact about concavification.

Fact A.1. *In the optimization problem,*

$$\text{cav } f(x) = \sup \left\{ \sum_m \lambda_m f(x_m) : \sum_m \lambda_m x_m = x \right\},$$

where f is defined on $X \subseteq \mathbb{R}^d$, the number of points can be restricted to $d + 1$. That is, without loss of generality, the supremum is taken over families $(\lambda_m, x_m)_{m=1}^{d+1}$.

The reader is referred to [Rockafellar \(1970, Corollary 17.1.5, p. 157\)](#). This implies that in a persuasion problem with unrestricted communication, the number of messages can be bounded by the dimension of $\Delta(\Omega)$ plus one, that is the number of states.

Corollary A.2. *In the optimisation problem,*

$$F^g(x, \gamma) = \sup \left\{ \sum_m \lambda_m f(x_m) : \sum_m \lambda_m x_m = x, \sum_m \lambda_m g(x_m) \geq \gamma \right\}$$

where f is defined on $X \subseteq \mathbb{R}^d$, the number of points can be restricted to $d + 2$.

This follows from [Corollary 3.4](#) and [Fact A.1](#), since the function f^g is defined on $X \times \mathbb{R} \subseteq \mathbb{R}^{d+1}$. Applying to the problem of optimal splitting under information constraint, gives a number of messages bounded by the dimension of $\Delta(\Omega)$ plus two, that is the number of states plus one. \square

A.2.4 Proof of [Corollary 3.5](#), upper bound $|A|$

Let $\tilde{A}(\nu) = \text{argmin} \left\{ \sum_{\omega} \nu(\omega) u_S(\omega, a) : a \in A^*(\nu) \right\}$ be the set of optimal actions of the receiver at ν which are worse for the sender.

Claim A.3. *For any action a , the set of ν 's such that $a \in \tilde{A}(\nu)$ is convex.*

Proof. Observe first that the set of ν 's such that $a \in A^*(\nu)$ is defined by linear inequalities, i.e. the optimality of a , therefore is convex. Consider now $a \in \tilde{A}(\nu_1) \cap \tilde{A}(\nu_2)$ and let us show that $a \in \tilde{A}(t\nu_1 + (1-t)\nu_2)$ for $t \in (0, 1)$. We have $a \in A^*(\nu_1) \cap A^*(\nu_2)$ and by the remark above, $a \in A^*(t\nu_1 + (1-t)\nu_2)$. Take $b \in A^*(t\nu_1 + (1-t)\nu_2)$. We thus have

$$\sum_{\omega} (t\nu_1(\omega) + (1-t)\nu_2(\omega)) u_R(\omega, a) = \sum_{\omega} (t\nu_1(\omega) + (1-t)\nu_2(\omega)) u_R(\omega, b).$$

Since $a \in A^*(\nu_1) \cap A^*(\nu_2)$,

$$\sum_{\omega} \nu_1(\omega) u_R(\omega, a) \geq \sum_{\omega} \nu_1(\omega) u_R(\omega, b), \quad \sum_{\omega} \nu_2(\omega) u_R(\omega, a) \geq \sum_{\omega} \nu_2(\omega) u_R(\omega, b).$$

Combined together, we get $b \in A^*(\nu_1) \cap A^*(\nu_2)$. Since $a \in \tilde{A}(\nu_1) \cap \tilde{A}(\nu_2)$,

$$\sum_{\omega} \nu_1(\omega) u_R(\omega, a) \leq \sum_{\omega} \nu_1(\omega) u_R(\omega, b), \quad \sum_{\omega} \nu_2(\omega) u_S(\omega, a) \leq \sum_{\omega} \nu_2(\omega) u_S(\omega, b).$$

Taking the convex combination of these two inequalities proves the claim. \square

Consider a feasible splitting (λ_m, μ_m) such that $\sum_m \lambda_m \nu_m = \mu$ and $\sum_m \lambda_m H(\nu_m) \geq H(\mu) - C$. For each action a , define $M(a) = \{m : \tilde{A}(\nu_m) = \{a\}\}$. Denote $\tilde{\lambda}_a = \sum_{m \in M(a)} \lambda_m$ and

$$\tilde{\nu}_a = \sum_{m \in M(a)} \frac{\lambda_m}{\tilde{\lambda}_a} \nu_m.$$

We have:

$$\begin{aligned} \mu &= \sum_m \lambda_m \nu_m \\ &= \sum_a \tilde{\lambda}_a \sum_{m \in M(a)} \frac{\lambda_m}{\tilde{\lambda}_a} \nu_m \\ &= \sum_a \tilde{\lambda}_a \tilde{\nu}_a. \end{aligned}$$

This defines a splitting of μ with $|A|$ elements. We argue that the payoff is the same as the initial splitting. Let us calculate the expected payoff. From the previous claim, for each action a , $a \in \tilde{A}(\tilde{\nu}_a)$. We thus have:

$$\begin{aligned} \sum_m \lambda_m u_S^*(\nu_m) &= \sum_a \tilde{\lambda}_a \sum_{m \in M(a)} \frac{\lambda_m}{\tilde{\lambda}_a} \sum_{\omega} \nu_m(\omega) u_S^*(\omega, a) \\ &= \sum_a \tilde{\lambda}_a \sum_{\omega} \tilde{\nu}_a(\omega) u_S^*(\omega, a) \\ &= \sum_a \tilde{\lambda}_a u_S^*(\tilde{\nu}_a). \end{aligned}$$

To conclude the proof, we check that the information constraint is satisfied. This follows from the concavity of entropy. Indeed,

$$H(\tilde{\nu}_a) \geq \sum_{m \in M(a)} \frac{\lambda_m}{\tilde{\lambda}_a} H(\nu_m)$$

and thus,

$$\sum_a \tilde{\lambda}_a H(\tilde{\nu}_a) \geq \sum_m \lambda_m H(\nu_m) \geq H(\mu) - C.$$

□

A.3 Proof of Theorem 3.1, point 1, the upper bound

1. For each pair of integers k, n , $U_S^*(\mu_n, Q_k) \leq V(\mu, \frac{k}{n}C(Q))$.

Proof. Let us fix a strategy σ of the sender. This induces a probability distribution \mathbb{P}_σ of sequences in $\Omega^n \times X^k \times Y^k$, the associated random sequences are denoted $(\boldsymbol{\omega}^n, \mathbf{x}^k, \mathbf{y}^k)$. Let \mathbf{t} be a uniformly distributed random variable over $\{1, \dots, n\}$, independent from $(\boldsymbol{\omega}^n, \mathbf{x}^k, \mathbf{y}^k)$ and denote $\mathbf{m} = (\mathbf{y}^k, \mathbf{t})$ taking values in $M = Y^k \times \{1, \dots, n\}$.

We denote $\tilde{\mathbb{P}}(\boldsymbol{\omega}, \mathbf{m})$ the joint probability distribution of $(\boldsymbol{\omega}, \mathbf{m})$ defined by:

$$\begin{aligned} \tilde{\mathbb{P}}(\boldsymbol{\omega}, \mathbf{m}) &= \tilde{\mathbb{P}}(\boldsymbol{\omega} = \omega, (\mathbf{y}^k, \mathbf{t}) = m) \\ &= \tilde{\mathbb{P}}(\mathbf{t} = t) \cdot \tilde{\mathbb{P}}(\boldsymbol{\omega} = \omega, \mathbf{y}^k = y^k | \mathbf{t} = t) \\ &= \frac{1}{n} \cdot \mathbb{P}_\sigma(\boldsymbol{\omega}_t = \omega, \mathbf{y}^k = y^k). \end{aligned}$$

Note that the marginal distribution of $\tilde{\mathbb{P}}(\boldsymbol{\omega}, \mathbf{m})$ on Ω is equal to the prior μ :

$$\begin{aligned} \tilde{\mathbb{P}}(\boldsymbol{\omega}) &= \sum_{t, y^k} \tilde{\mathbb{P}}(\boldsymbol{\omega} = \omega, \mathbf{y}^k = y^k, \mathbf{t} = t) \\ &= \sum_{t, y^k} \frac{1}{n} \cdot \mathbb{P}_\sigma(\boldsymbol{\omega}_t = \omega, \mathbf{y}^k = y^k) \\ &= \sum_{t=1}^n \frac{1}{n} \cdot \mathbb{P}_\sigma(\boldsymbol{\omega}_t = \omega) \\ &= \mathbb{P}_\sigma(\boldsymbol{\omega}) \cdot \sum_{t=1}^n \frac{1}{n} = \mu(\boldsymbol{\omega}). \end{aligned}$$

Fix now a strategy τ of the receiver $\tau : Y^k \rightarrow A^n$ and define $\tilde{\tau} : M \rightarrow A$ where $\tilde{\tau}(m) = \tilde{\tau}(\mathbf{y}^k, t) =$

$\tau_t(y^k)$, the t -th coordinate of $\tau(y^k)$. The expected average payoff of player $i = R, S$ writes:

$$\mathbb{E}_{\sigma, \tau}[\bar{u}_i] = \sum_{\omega^n, x^k, y^k} \mathbb{P}_{\sigma}(\omega^n, x^k, y^k) \left[\frac{1}{n} \sum_{t=1}^n u_i(\omega_t, \tau_t(y^k)) \right] \quad (5)$$

$$= \sum_{t=1}^n \sum_{\omega_t, x^k, y^k} \frac{1}{n} \cdot \mathbb{P}_{\sigma}(\omega_t, x^k, y^k) \cdot u_i(\omega_t, \tau_t(y^k)) \quad (6)$$

$$= \sum_{t=1}^n \sum_{\omega_t, y^k} \frac{1}{n} \cdot \mathbb{P}_{\sigma}(\omega_t, y^k) \cdot u_i(\omega_t, \tau_t(y^k)) \quad (7)$$

$$= \sum_{\omega, y^k, t} \tilde{\mathbb{P}}(\omega, y^k, t) \cdot u_i(\omega, \tilde{\tau}(y^k, t)) \quad (8)$$

$$= \sum_{\omega, m} \tilde{\mathbb{P}}(\omega, m) \cdot u_i(\omega, \tilde{\tau}(m)). \quad (9)$$

Equation (6) implies Equation (7) by summing over x^k which does not enter the payoff function. All other steps are reorderings and change of variables.

A strategy τ is a best-reply to σ if and only if:

$$\begin{aligned} \tau(y^k) &\in \arg \max_{a^n \in A^n} \sum_{\omega^n, x^k, y^k} \mu(\omega^n) \sigma(x^k | \omega^n) Q(y^k | x^k) \bar{u}_R(\omega^n, a^n) \\ \iff \tilde{\tau}(m) &\in \arg \max_{a \in A} \sum_{\omega, m} \tilde{\mathbb{P}}(\omega, m) \cdot u_R(\omega, a) \\ \iff \tilde{\tau}(m) &\in \arg \max_{a \in A} \sum_{\omega} \tilde{\nu}_{\sigma}(\omega | m) \cdot u_R(\omega, a) \\ \iff \tilde{\tau}(m) &\in A^*(\tilde{\nu}_{\sigma}(\cdot | m)) \end{aligned}$$

where $\tilde{\nu}_{\sigma}(\omega | m) = \tilde{\mathbb{P}}(\omega | m)$. We deduce for any strategy σ of the sender and any best-reply τ of the sender, the expected average payoffs are those induced by the splitting:

$$\mu(\omega) = \sum_m \tilde{\mathbb{P}}(m) \tilde{\nu}_{\sigma}(\omega | m).$$

Now, we bound the mutual information of this splitting. Throughout the proof, we will abuse our notations in a way that is common in information theory (see [Cover and Thomas \(2006\)](#)). When \mathbf{x} is a random variable with distribution $p(x)$, we write $H(\mathbf{x})$ for $H(p)$, when (\mathbf{x}, \mathbf{y}) is a pair of random variables with joint distribution $p(x, y)$, we write $H(\mathbf{y} | \mathbf{x})$ for $\sum_x p(x) H(p(\cdot | x))$. Last, we will write $I(\mathbf{x}; \mathbf{y})$ without explicit reference to the joint distribution.

For any strategy σ , we have:

$$0 \leq I(\mathbf{x}^k; \mathbf{y}^k) - I(\boldsymbol{\omega}^n; \mathbf{y}^k) \quad (10)$$

$$= \sum_{t=1}^k H(\mathbf{y}_t | \mathbf{y}^{t-1}) - \sum_{t=1}^k H(\mathbf{y}_t | \mathbf{x}^k, \mathbf{y}^{t-1}) - \sum_{t=1}^n H(\boldsymbol{\omega}_t | \boldsymbol{\omega}^{t-1}) + \sum_{t=1}^n H(\boldsymbol{\omega}_t | \mathbf{y}^k, \boldsymbol{\omega}^{t-1}) \quad (11)$$

$$\leq \sum_{t=1}^k H(\mathbf{y}_t) - \sum_{t=1}^k H(\mathbf{y}_t | \mathbf{x}_t) - n \cdot H(\boldsymbol{\omega}) + \sum_{t=1}^n H(\boldsymbol{\omega}_t | \mathbf{y}^k) \quad (12)$$

$$= \sum_{t=1}^k I(\mathbf{x}_t; \mathbf{y}_t) - n \cdot H(\boldsymbol{\omega}) + n \cdot \sum_{t=1}^n \mathbb{P}(\mathbf{t} = t) \cdot H(\boldsymbol{\omega} | \mathbf{y}^k, \mathbf{t} = t) \quad (13)$$

$$\leq k \cdot \max_{\mathbb{P}(x)} I(\mathbf{x}; \mathbf{y}) - n \cdot H(\boldsymbol{\omega}) + n \cdot H(\boldsymbol{\omega} | \mathbf{y}^k, \mathbf{t}) \quad (14)$$

$$= k \cdot \max_{\mathbb{P}(x)} I(\mathbf{x}; \mathbf{y}) - n \cdot H(\boldsymbol{\omega}) + n \cdot H(\boldsymbol{\omega} | \mathbf{m}) \quad (15)$$

$$= k \cdot \max_{\mathbb{P}(x)} I(\mathbf{x}; \mathbf{y}) - n \cdot I(\boldsymbol{\omega}; \mathbf{m}). \quad (16)$$

- Equation (10) holds since the triple $(\boldsymbol{\omega}^n, \mathbf{x}^k, \mathbf{y}^k)$ has the Markov chain property; that is, its joint distribution writes $\mu(\boldsymbol{\omega}^n) \sigma(x^k | \boldsymbol{\omega}^n) Q(\mathbf{y}^k | x^k)$. This implies $I(\mathbf{x}^k; \mathbf{y}^k) \geq I(\boldsymbol{\omega}^n; \mathbf{y}^k)$, that is \mathbf{x}^k is more informative than $\boldsymbol{\omega}^n$ about \mathbf{y}^k (Cover and Thomas, 2006, Theorem 2.8.1, p. 34).

- Equation (11) comes from the chain rule of entropy $H(\mathbf{y}^k) = \sum_{t=1}^k H(\mathbf{y}_t | \mathbf{y}^{t-1})$.

- Equation (12) follows since the channel is memoryless $H(\mathbf{y}_t | \mathbf{x}^k, \mathbf{y}^{t-1}) = H(\mathbf{y}_t | \mathbf{x}_t)$, the sequence of states is i.i.d. $H(\boldsymbol{\omega}_t | \boldsymbol{\omega}^{t-1}) = H(\boldsymbol{\omega}_t)$, and conditioning reduces entropy $H(\boldsymbol{\omega}_t | \mathbf{y}^k, \boldsymbol{\omega}^{t-1}) \leq H(\boldsymbol{\omega}_t | \mathbf{y}^k)$.

- Equation (13) is a simple rewriting with the introduction of the uniform random variable $\mathbf{t} \in \{1, \dots, n\}$.

- Equation (14) comes from taking the maximum over the marginal distribution $\mathbb{P}(x)$.

- Equation (15) comes from the change of variable $\mathbf{m} = (\mathbf{y}^k, \mathbf{t})$.

Then, Equation (16) is equivalent to:

$$\begin{aligned} & k \cdot \max_{\mathbb{P}(x)} I(\mathbf{x}; \mathbf{y}) - n \cdot I(\boldsymbol{\omega}; \mathbf{m}) \geq 0 \\ \iff & H(\boldsymbol{\omega} | \mathbf{m}) \geq H(\boldsymbol{\omega}) - \frac{k}{n} \cdot \max_{\mathbb{P}(x)} I(\mathbf{x}; \mathbf{y}) \\ \iff & \sum_m \lambda_m H(\mu_m) \geq H(\mu) - \frac{k}{n} \cdot C(Q). \end{aligned}$$

Therefore, for any strategy σ and all n, k , we have:

$$\begin{aligned}
& \min_{\tau \in BR(\sigma)} \sum_{\omega^n, x^k, y^k} \mu(\omega^n) \sigma(x^k | \omega^n) Q(y^k | x^k) \bar{u}_S(\omega^n, \tau(y^k)) \\
&= \min_{\tilde{\tau} \in BR(\sigma)} \sum_{\omega, m} \tilde{\mathbb{P}}(\omega, m) \cdot u_S(\omega, \tilde{\tau}(m)) \\
&= \sum_m \tilde{\mathbb{P}}(m) \min_{\tilde{\tau}(m) \in A^*(\tilde{\nu}_\sigma(\cdot | m))} \sum_{\omega} \tilde{\nu}_\sigma(\cdot | m) \cdot u_S(\omega, \tilde{\tau}(m)) \\
&= \sum_m \tilde{\mathbb{P}}(m) \cdot u_S^*(\tilde{\nu}_\sigma(\cdot | m)) \\
&\leq \sup_{\sigma} \left\{ \sum_m \tilde{\mathbb{P}}(m) \cdot u_S^*(\tilde{\nu}_\sigma(\cdot | m)) \right. \\
&\quad \text{s.t. } \sum_m \tilde{\mathbb{P}}(m) \cdot \tilde{\nu}_\sigma(\cdot | m) = \mu, \\
&\quad \left. \text{and } \sum_m \tilde{\mathbb{P}}(m) \cdot H(\tilde{\nu}_\sigma(\cdot | m)) \geq H(\mu) - \frac{k}{n} \cdot C(Q) \right\} \\
&= \sup \left\{ \sum_m \lambda_m \cdot u_S^*(\nu_m) \right. \\
&\quad \text{s.t. } \sum_m \lambda_m \nu_m = \mu, \\
&\quad \left. \text{and } \sum_m \lambda_m H(\nu_m) \geq H(\mu) - \frac{k}{n} \cdot C(Q) \right\} \\
&= V(\mu, \frac{k}{n} C(Q)).
\end{aligned}$$

This proves that for all n and k we have:

$$\begin{aligned}
U_S^*(\mu_n, Q_k) &= \sup_{\sigma} \min_{\tau \in BR(\sigma)} \sum_{\omega^n, x^k, y^k} \mu(\omega^n) \sigma(x^k | \omega^n) Q(y^k | x^k) \bar{u}_S(\omega^n, \tau(y^k)) \\
&\leq V(\mu, \frac{k}{n} C(Q))
\end{aligned}$$

as desired. □

A.4 Proof of Theorem 3.1, point 2, the limit value

2. For each $r \in [0, +\infty]$ and each pair of sequences $(k_j, n_j)_{j \in \mathbb{N}}$ such that $\lim_{j \rightarrow \infty} \max(n_j, k_j) = \infty$ and $\lim_{j \rightarrow \infty} \frac{k_j}{n_j} = r$, we have $\lim_{j \rightarrow \infty} U_S^*(\mu^{n_j}, Q^{k_j}) = V(\mu, rC(Q))$.

For this proof we will consider the case r finite, $n_j \rightarrow \infty$, $k_j \rightarrow \infty$, $\frac{k_j}{n_j} \rightarrow r$. The proof for the case where $n_j \rightarrow \infty$, $k_j \rightarrow \infty$, $\frac{k_j}{n_j} \rightarrow \infty$ is a consequence by considering r finite large enough such that $rC(Q) > H(\mu)$. For the cases where either n_j or k_j is bounded, see Section 3.1.2. In the rest of the proof, we will consider pairs of integers (k, n) which are generic terms of such a

sequence $(k_j, n_j)_{j \in \mathbb{N}}$ and omit the index j for simplicity of notations.

A.4.1 Zero capacity.

First, we investigate the case $C(Q) = 0$.

Lemma A.4. *If the channel capacity is equal to zero $\max_{p(x)} I(\mathbf{x}; \mathbf{y}) = 0$, then for all k, n , we have:*

$$U_S^*(\mu^n, Q^k) = V(\mu, \frac{k}{n}C(Q)) = u_S^*(\mu).$$

Proof. [Lemma A.4] Let (\mathbf{x}, \mathbf{y}) be a pair of random variables such that the conditional probability of $\{\mathbf{y} = y\}$ given $\{\mathbf{x} = x\}$ is $Q(y|x)$. If the capacity of the channel is 0, then $I(\mathbf{x}; \mathbf{y}) = H(\mathbf{y}) - H(\mathbf{y}|\mathbf{x}) = 0$ which implies that \mathbf{x} and \mathbf{y} are independent: no information can be sent through the channel. This implies that for any splitting which satisfies the information constraint, the random variables $\boldsymbol{\omega}$ and \mathbf{m} are independent, and for all $m \in M$ we have $\nu_m = \mu$. Hence:

$$V(\mu, \frac{k}{n}C(Q)) = u_S^*(\mu).$$

Moreover, for any strategy σ , the sequence of messages \mathbf{y}^k of the receiver is independent from the sequence of states $\boldsymbol{\omega}^n$. It follows that:

$$\begin{aligned} U_S^*(\mu^n, Q^k) &= \sup_{\sigma} \min_{\tau \in BR(\sigma)} \sum_{\omega^n, x^k, y^k} \mu^n(\omega^n) \sigma(x^k | \omega^n) Q^k(y^k) \bar{u}_S(\omega^n, \tau(y^k)) \\ &= \min_{\tau \in BR(\sigma)} \sum_{\omega^n, y^k} \mu^n(\omega^n) Q^k(y^k) \left[\frac{1}{n} \sum_{t=1}^n u_S(\omega_t, \tau_t(y^k)) \right] \\ &= \frac{1}{n} \sum_{t=1}^n \min_{a_t \in A^*(\mu)} \sum_{\omega_t} \mu(\omega_t) u_S(\omega_t, a_t) \\ &= \min_{a \in A^*(\mu)} \sum_{\omega} \mu(\omega) u_S(\omega, a) = u_S^*(\mu), \end{aligned}$$

which concludes the proof. □

A.4.2 Positive channel capacity.

We assume from now on $C(Q) > 0$. The goal is to take a splitting of the prior which satisfies the information constraint and to show that the associated payoff can be approximately achieved by strategy σ of the sender and a best-reply $\tau \in BR(\sigma)$ of the receiver. The next lemma

states that we can focus on splittings such that the information constraint is satisfied with strict inequality and where the action of the receiver is unique for each posterior. Concretely, we prove that such splittings are dense in the set of feasible splittings. Recall that we denote $\tilde{A}(\nu)$ the set of worst optimal actions when the belief is $\nu \in \Delta(\Omega)$:

$$\tilde{A}(\nu) = \operatorname{argmin} \left\{ \sum_{\omega} \nu(\omega) u_S(\omega, a) : a \in A^*(\nu) \right\}.$$

Consider the following program:

$$\begin{aligned} \hat{V}(\mu, \frac{k}{n}C(Q)) = \sup & \left\{ \sum_m \lambda_m u_S^*(\nu_m) \right. \\ & \text{s.t. } \sum_m \lambda_m \nu_m = \mu, \\ & \text{and } H(\mu) - \sum_m \lambda_m H(\nu_m) < \frac{k}{n}C(Q) \\ & \left. \text{and } \forall m, \tilde{A}(\nu_m) \text{ is a singleton} \right\}. \end{aligned}$$

Lemma A.5. *For all integers (k, n) , $\mu \in \Delta(\Omega)$ and Q such that $C(Q) > 0$ we have:*

$$V(\mu, \frac{k}{n}C(Q)) = \hat{V}(\mu, \frac{k}{n}C(Q)). \quad (17)$$

The proof of Lemma A.5 is postponed to Section A.4.3. Then, the proof of our main result continues with two lemmas. In Lemma A.8, we approximate the payoff yielded by any strategy. We will see the relevance of the number of stages where the actual belief of the receiver is close to the desired target, and the importance of this number being large. Next, in Lemma A.9 we prove that there is a strategy for which this holds. We use there known results from information theory for defining the coding scheme which gives the strategy of the sender. Then, we prove that this strategy actually controls the Bayesian beliefs of the receiver.

Given a strategy σ of the sender, we denote the induced expected payoff as follows:

$$\begin{aligned} \hat{U}_{S,\sigma}(\mu^n, Q^k) &= \min_{\tau \in BR(\sigma)} \sum_{\omega^n, x^k, y^k} \mu(\omega^n) \sigma(x^k | \omega^n) Q(y^k | x^k) \bar{u}_S(\omega^n, \tau(y^k)), \\ &= \min_{\tau \in BR(\sigma)} \mathbb{E}_{\omega^n, x^k, y^k} [\bar{u}_S(\omega^n, \tau(y^k))]. \end{aligned}$$

Let $\nu_{t,y^k}^\sigma \in \Delta(\Omega)$ denote the posterior belief on ω_t conditional on the sequence y^k . That is,

$$\nu_{t,y^k}^\sigma(\omega) = \mathbb{P}_\sigma(\omega_t = \omega | y^k).$$

For $\nu_1, \nu_2 \in \Delta(\Omega)$, the Kullback-Leibler (KL) divergence is,

$$D(\nu_1 \parallel \nu_2) = \sum_{\omega} \nu_1(\omega) \log \frac{\nu_1(\omega)}{\nu_2(\omega)}.$$

We will introduce several positive parameters α, γ, δ , to be thought of as small.

Notation A.6. For a sequence (m^n, y^k) and $\alpha > 0$, denote

$$T_{\alpha}(m^n, y^k) = \left\{ t \in \{1, \dots, n\} : D(\nu_{t, y^k}^{\sigma} \parallel \nu_{m_t}) \leq \frac{\alpha^2}{2 \ln 2} \right\}.$$

This is the set of indices $t = 1, \dots, n$ such that the posterior belief ν_{t, y^k}^{σ} about ω_t is close to the theoretical belief ν_{m_t} . Intuitively, this is the set of indices where *the message* m_t is approximately transmitted. Now, we define an event $B_{\alpha, \gamma, \delta} \subseteq M^n \times Y^k$ such that for every $(m^n, y^k) \in B_{\alpha, \gamma, \delta}$, $\frac{1}{n} \sum_{t=1}^n u_S^*(\nu_{t, y^k}^{\sigma})$ is close to $\sum_m \lambda_m u_S^*(\nu_m)$.

Notation A.7. For a sequence (m^n, y^k) and $m \in M$, denote

$$\text{freq}_m(m^n, y^k) = \frac{1}{n} \left| \{t = 1, \dots, n : m_t = m\} \right|$$

the empirical frequency of message m in the sequence m^n . For $\alpha, \gamma, \delta > 0$, let

$$B_{\alpha, \gamma, \delta} = \left\{ (m^n, y^k) : \frac{|T_{\alpha}(m^n, y^k)|}{n} \geq 1 - \gamma \text{ and } \sum_m |\lambda_m - \text{freq}_m(m^n, y^k)| \leq \delta \right\}$$

Lemma A.8.

$$\left| \widehat{U}_{S, \sigma}(\mu^n, Q^k) - \widehat{V}(\mu, rC(Q)) \right| \leq (\alpha + 2\gamma + \delta) \|u\| + (1 - \mathbb{P}_{\sigma}(B_{\alpha, \gamma, \delta})) \|u\|.$$

The proof of Lemma A.8 is given in Section A.4.4. We see from this inequality that estimating the probability of the set $B_{\alpha, \gamma, \delta}$ is crucial and that we would like the probability of the complement $\mathbb{P}_{\sigma}(B_{\alpha, \gamma, \delta}^c)$ to be small.

Last, Lemma A.9 corresponds to the actual construction of the strategy.

Lemma A.9. *Assume that the splitting $(\lambda_m, \nu_m)_m$ satisfies the three conditions:*

$$\sum_m \lambda_m \nu_m = \mu, \quad (18)$$

$$H(\mu) - \sum_m \lambda_m H(\mu_m) < rC(Q), \quad (19)$$

$$\forall m, \tilde{A}(\nu_m) \text{ is a singleton}, \quad (20)$$

then $\forall \varepsilon > 0, \forall \alpha > 0, \forall \gamma > 0, \exists \bar{\delta}, \forall \delta \leq \bar{\delta}, \exists \bar{n}, \forall n \geq \bar{n}, \exists \sigma$, such that $\mathbb{P}_\sigma(B_{\alpha, \gamma, \delta}^c) \leq \varepsilon$.

The proof of Lemma A.9 is in Appendix A.4.5. The idea is that, since the information constraint is satisfied i.e. $I(\omega; \mathbf{m}) < rC(Q)$, there is enough capacity to transmit \mathbf{m} over the channel. More precisely, we construct a strategy such that the set $B_{\alpha, \gamma, \delta}$ has probability close to 1. This way, for most sequences $(\omega^n, m^n, x^k, y^k)$, the receiver gets the right message in most stages. That is, at most stages the receiver plays the action corresponding to the message.

We may now conclude the main proof. We combine the inequality of Lemma A.8 with the bound $\mathbb{P}_\sigma(B_{\alpha, \gamma, \delta}^c) \leq \varepsilon$ of Lemma A.9. We choose the parameters $\alpha, \gamma, \eta, \delta$ small and then n large in order to obtain the following:

Proposition A.10. *For all $r > 0$ and $\varepsilon > 0$, there exists integers $N(\varepsilon), K(\varepsilon)$ such that for all $n \geq N(\varepsilon)$, $k \geq K(\varepsilon)$ and $|\frac{k}{n} - r| \leq \varepsilon$, there exists a strategy σ such that:*

$$\left| \widehat{U}_{S, \sigma}(\mu^n, Q^k) - \widehat{V}(\mu, rC(Q)) \right| \leq \varepsilon. \quad (21)$$

With Lemma A.5, this ends the proof of point 2 of Theorem 3.1. □

A.4.3 Proof of Lemma A.5

Remark A.11. From Corollary 3.5, we know that we can restrict the number of messages, i.e. the number of posteriors to $K = \min\{|A|, |\Omega| + 1\}$. Therefore, from now on a splitting $(\lambda_m, \nu_m)_m$ will be understood to be a composed of $\lambda = (\lambda_1, \dots, \lambda_K) \in \Delta(\{1, \dots, K\})$ and $(\nu_m)_m \in (\Delta(\Omega))^K$. The set of splittings of μ is thus a convex and compact subset of $\Delta(\{1, \dots, K\}) \times (\Delta(\Omega))^K$ which itself is a compact and convex set in some finite dimension space. All statements below about closed or open sets of splittings relate to the topology induced by the Euclidean topology on this finite dimension space.

We consider the following sets:

$$\begin{aligned}
\mathcal{S}_1 &= \left\{ (\lambda_m, \nu_m)_m, \quad \text{s.t.} \quad \sum_m \lambda_m \nu_m = \mu, \right. \\
&\quad \left. \text{and} \quad \sum_m \lambda_m H(\nu_m) \geq H(\mu) - \frac{k}{n} C(Q) \right\}, \\
\mathcal{S}_2 &= \left\{ (\lambda_m, \nu_m)_m, \quad \text{s.t.} \quad \sum_m \lambda_m \nu_m = \mu, \right. \\
&\quad \left. \text{and} \quad \forall m, \tilde{A}(\nu_m) \text{ is a singleton} \right\}, \\
\mathcal{S}_3 &= \left\{ (\lambda_m, \nu_m)_m, \quad \text{s.t.} \quad \sum_m \lambda_m \nu_m = \mu, \right. \\
&\quad \left. \text{and} \quad \sum_m \lambda_m H(\nu_m) > H(\mu) - \frac{k}{n} C(Q) \right\}.
\end{aligned}$$

We will prove that the set $\mathcal{S}_2 \cap \mathcal{S}_3$ is dense in \mathcal{S}_1 , which will imply that Equation (17) is satisfied. We first argue that $\tilde{A}(\nu)$ is a singleton for an open and dense set of posteriors ν .

Definition A.12. *Two actions a and b are equivalent $a \sim_i b$ for player $i = S, R$, if for all $\omega \in \Omega$, $u_i(\omega, a) = u_i(\omega, b)$.*

We say that two actions a and b are *completely equivalent* if they are equivalent for both players. Without loss of generality, we assume that no two actions are completely equivalent. Otherwise, we can merge them into one single action and work on the reduced problem.

Denote $F_i \subseteq \Delta(\Omega)$ the set of beliefs for which player $i \in \{S, R\}$ is indifferent between two actions which are not equivalent:

$$F_i = \left\{ \nu \in \Delta(\Omega) : \exists a, b, a \not\sim_i b, \sum_{\omega} \nu(\omega) u_i(\omega, a) = \sum_{\omega} \nu(\omega) u_i(\omega, b) \right\}.$$

Let $F^c = \Delta(\Omega) \setminus (F_R \cup F_S)$ be the set of beliefs where at least one player is not indifferent between any two actions.

Claim A.13. *The set F^c is open and dense in $\Delta(\Omega)$ and for each $\nu \in F^c$, $\tilde{A}(\nu)$ is a singleton.*

Proof. [Claim A.13] For each i and each pair of actions a, b with $a \not\sim_i b$, the set,

$$F_i(a, b) = \left\{ \nu \in \Delta(\Omega) : \sum_{\omega} \nu(\omega) u_i(\omega, a) = \sum_{\omega} \nu(\omega) u_i(\omega, b) \right\}$$

is a closed hyperplane of dimension $\dim(F_i(a, b)) \leq |\Omega| - 2$. Thus, F_R and F_S are closed and $F_R \cup F_S$ is included in a finite union of hyperplanes of dimension at most $|\Omega| - 2$. The

complementary set is thus open and dense in $\Delta(\Omega)$.

Then, if $\tilde{A}(\nu)$ contains two distinct actions $a \neq b$, both players are indifferent between a and b at ν . Thus, if $\nu \in F^c$, $\tilde{A}(\nu)$ is a singleton. \square

It follows that \mathcal{S}_2 is open and dense in \mathcal{S}_1 .

Claim A.14. *If the channel capacity is strictly positive $C(Q) > 0$, the set \mathcal{S}_3 is nonempty, open and dense in \mathcal{S}_1 .*

Proof. [Claim A.14] Take a feasible splitting $(\lambda_m, \nu_m)_m$ in \mathcal{S}_1 :

$$\sum_m \lambda_m H(\nu_m) \geq H(\mu) - \frac{k}{n} C(Q).$$

For $\varepsilon > 0$, consider the perturbed splitting $(\lambda_m, (1-\varepsilon)\nu_m + \varepsilon\mu)_m$. From concavity of the entropy,

$$\begin{aligned} \sum_m \lambda_m H((1-\varepsilon)\nu_m + \varepsilon\mu) &\geq (1-\varepsilon) \sum_m \lambda_m H(\nu_m) + \varepsilon H(\mu), \\ &\geq H(\mu) - \frac{k}{n} C(Q) + \varepsilon \frac{k}{n} C(Q) \\ &> H(\mu) - \frac{k}{n} C(Q); \end{aligned}$$

thus, the information constraint is satisfied with strict inequality for $\varepsilon > 0$. It follows that \mathcal{S}_3 is nonempty and dense in \mathcal{S}_1 . By continuity of the entropy, \mathcal{S}_3 is open in \mathcal{S}_1 . \square

Since \mathcal{S}_2 and \mathcal{S}_3 are open and dense, $\mathcal{S}_2 \cap \mathcal{S}_3$ is also open and dense in \mathcal{S}_1 . We can conclude that $V(\mu, \frac{k}{n}C(Q)) = \hat{V}(\mu, \frac{k}{n}C(Q))$ as desired. This follows from the fact that the function

$$u_S^*(\nu) = \min_{a \in A^*(\nu)} \sum_{\omega} \nu(\omega) u_S(\omega, a)$$

is lower-semi continuous and the supremum of an l.s.c. function over a dense set is the supremum over the full set.

It should be noticed that this is the only argument in the proof where the assumption that the receiver chooses the worst action for the sender, has a bite. When the receiver chooses the best action for the sender, we should consider $u_S^{**}(\nu) = \max_{a \in A^*(\nu)} \sum_{\omega} \nu(\omega) u_S(\omega, a)$ which is upper-semi continuous. In that case, the supremum over the dense set $\mathcal{S}_2 \cap \mathcal{S}_3$ might be less than the supremum over \mathcal{S}_1 . However, this can only happen when the information constraint is binding at optimum and all posteriors in the optimal splitting are points of indifference for the receiver. This case is nongeneric in our class of persuasion problems: a slight change of the

payoff function of the receiver would perturb the points of indifference and thus the points of discontinuity of u^* and u^{**} . \square

A.4.4 Proof of Lemma A.8

The strategy σ induces a joint probability distribution \mathbb{P}_σ over $\Omega^n \times M^n \times X^k \times Y^k$:

$$\mathbb{P}_\sigma(\omega^n, m^n, x^k, y^k) = \prod_{t=1}^n \mu(\omega_t) \times \sigma(m^n, x^k | \omega^n) \times \prod_{t=1}^n Q(y_t | x_t).$$

For each sequence y^k of messages and for each t , the receiver chooses an optimal action $a_t \in A^*(\nu_{t,y^k}^\sigma)$. In the worst case (for the sender), this action a_t belongs to $\tilde{A}(\nu_{t,y^k}^\sigma)$. It follows that:

Claim A.15.

$$\hat{U}_{S,\sigma}(\mu^n, Q^k) = \sum_{m^n, y^k} \mathbb{P}_\sigma(m^n, y^k) \frac{1}{n} \sum_{t=1}^n u_S^*(\nu_{t,y^k}^\sigma).$$

Remark A.16. *Since the set of posteriors ν such that $\tilde{A}(\nu)$ is a singleton is open, there exists $\alpha_0 > 0$ such that for all m :*

$$D(\nu \| \nu_m) \leq \alpha_0 \Rightarrow \tilde{A}(\nu) = \tilde{A}(\nu_m).$$

Whenever $\tilde{A}(\nu)$ is a singleton, denote $\tilde{A}(\nu) = \{\tilde{a}(\nu)\}$ the unique (worst) optimal action. From now on, we assume that $\alpha \in (0, \alpha_0)$. With the remark above, this implies that for each $t \in T_\alpha(m^n, y^k)$, the action chosen by the receiver for problem t is $\tau_t(m^n, y^k) = \tilde{a}(\nu_{m_t})$. So precisely, $T_\alpha(m^n, y^k)$ is the set of indices t such that the receiver plays the action $\tilde{a}(\nu_{m_t})$ which corresponds to the message m_t . In this sense, this is the set of indices for which the information transmission is successful.

Lemma A.17. *For each $(m^n, y^k) \in B_{\alpha,\gamma,\delta}$,*

$$\left| \frac{1}{n} \sum_{t=1}^n u_S^*(\nu_{t,y^k}^\sigma) - \sum_m \lambda_m u_S^*(\nu_m) \right| \leq (\alpha + 2\gamma + \delta) \|u\|,$$

where $\|u\| = \max_{\omega,a} |u_S(\omega, a)|$ is the largest absolute value of payoffs for the sender.

Proof. Denote $u^* = \sum_m \lambda_m u_S^*(\nu_m)$. We have:

$$\begin{aligned} \left| \frac{1}{n} \sum_{t=1}^n u_S^*(\nu_{t,y^k}^\sigma) - u^* \right| &\leq \left| \frac{1}{n} \sum_{t \in T_\alpha(m^n, y^k)} (u_S^*(\nu_{t,y^k}^\sigma) - u^*) \right| + \left| \frac{1}{n} \sum_{t \notin T_\alpha(m^n, y^k)} (u_S^*(\nu_{t,y^k}^\sigma) - u^*) \right| \\ &\leq \left| \frac{1}{n} \sum_{t \in T_\alpha(m^n, y^k)} (u_S^*(\nu_{t,y^k}^\sigma) - u^*) \right| + \gamma \|U\| \end{aligned}$$

Then:

$$\begin{aligned} \left| \frac{1}{n} \sum_{t \in T_\alpha(m^n, y^k)} (u_S^*(\nu_{t,y^k}^\sigma) - u^*) \right| &\leq \left| \frac{1}{n} \sum_{t \in T_\alpha(m^n, y^k)} (u_S^*(\nu_{t,y^k}^\sigma) - u_S^*(\nu_{m_t})) \right| \\ &\quad + \left| \frac{1}{n} \sum_{t \in T_\alpha(m^n, y^k)} (u_S^*(\nu_{m_t}) - u^*) \right| \end{aligned}$$

Since $\alpha \leq \alpha_0$, for each $t \in T_\alpha(m^n, y^k)$, $\tilde{a}(\nu_{t,y^k}^\sigma) = \tilde{a}(\nu_{m_t})$. Therefore, for $t \in T_\alpha(m^n, y^k)$

$$\left| u_S^*(\nu_{t,y^k}^\sigma) - u_S^*(\nu_{m_t}) \right| \leq \sum_\omega |\nu_{t,y^k}^\sigma(\omega) - \nu_{m_t}(\omega)| \cdot |u_S(\omega, a)| \leq \|\nu_{t,y^k}^\sigma - \nu_{m_t}\| \cdot \|u\| \leq \alpha \|u\|,$$

where the latter inequality comes from Pinsker's inequality⁹: $\|\nu_1 - \nu_2\| \leq \sqrt{2 \ln 2 D(\nu_1 \| \nu_2)}$ and the definition of $T_\alpha(m^n, y^k)$. It follows:

$$\left| \frac{1}{n} \sum_{t \in T_\alpha(m^n, y^k)} (u_S^*(\nu_{t,y^k}^\sigma) - u^*) \right| \leq \alpha \|u\| + \left| \frac{1}{n} \sum_{t \in T_\alpha(m^n, y^k)} (u_S^*(\nu_{m_t}) - u^*) \right|$$

Now from $\frac{|T_\alpha(m^n, y^k)|}{n} \geq 1 - \gamma$, we have:

$$\left| \frac{1}{n} \sum_{t \in T_\alpha(m^n, y^k)} (u_S^*(\nu_{m_t}) - u^*) \right| \leq \left| \frac{1}{n} \sum_{t=1}^n (u_S^*(\nu_{m_t}) - u^*) \right| + \gamma \|u\|.$$

Then:

$$\begin{aligned} \left| \frac{1}{n} \sum_{t=1}^n (u_S^*(\nu_{m_t}) - u^*) \right| &= \left| \sum_m (\text{freq}_m(m^n, y^k) - \lambda_m) u_S^*(\nu_m) \right| \\ &\leq \sum_m \left| \text{freq}_m(m^n, y^k) - \lambda_m \right| \cdot \left| u_S^*(\nu_m) \right| \\ &\leq \|u\| \delta. \end{aligned}$$

Collecting all inequalities together yields the desired conclusion. \square

⁹Cover and Thomas, 2006, Lemma 11.6.1, p. 370.

A.4.5 Proof of Lemma A.9

By hypothesis, the splitting $(\lambda_m, \nu_m)_m$ satisfies the three conditions:

$$\sum_m \lambda_m \nu_m = \mu, \quad (22)$$

$$H(\mu) - \sum_m \lambda_m H(\mu_m) < rC(Q), \quad (23)$$

$$\forall m, \tilde{A}(\nu_m) \text{ is a singleton.} \quad (24)$$

Let $M = \{1, \dots, |M|\}$ be the set of messages associated with this splitting.

Part 1. Coding scheme. We turn now to the actual construction. We use standard information theoretic techniques for Channel Coding (Gamal and Kim, 2011, Chap. 3.1, p. 38) and Lossy Source Coding (Gamal and Kim, 2011, Chap. 3.6, p. 56). Using information theoretic language, the sender is viewed as an *encoder* who encrypts his intended m^n messages in sequences of inputs x^k . The messages m^n are immaterial and can be seen as a pure mental construct of the sender. The encoding is such that a *decoder* who reads the sequence y^k , is able to determine the correct m^n with high probability. This is described as follows.

For $\delta > 0$, we define the set of *typical sequences* A_δ as follows:

$$A_\delta = \left\{ (\omega^n, m^n, x^k, y^k), \quad \text{s.t.} \quad \sum_{\omega, m} \left| \lambda_m \mu_m(\omega) - \text{freq}_{\omega, m}(\omega^n, m^n) \right| \leq \delta, \quad (25) \right.$$

$$\left. \text{and} \quad \sum_{x, y} \left| \mathbb{P}(x) \times Q(y|x) - \text{freq}_{x, y}(x^k, y^k) \right| \leq \delta \right\}. \quad (26)$$

A pair of sequences (ω^n, m^n) which satisfies Equation (25) will be called *jointly typical*. Similarly, pair of sequences (x^k, y^k) which satisfies Equation (26) will be called *jointly typical*. With a slight abuse of notation, we will write $(\omega^n, m^n) \in A_\delta$ or $(x^k, y^k) \in A_\delta$ to indicate jointly typical sequences.

Since condition (23) is satisfied with strict inequality, there exists a small parameter $\eta > 0$ and a “rate” $R \geq 0$, such that:

$$R = H(\mu) - \sum_m \lambda_m H(\mu_m) + \eta, \quad (27)$$

$$R \leq rC(Q) - \eta. \quad (28)$$

Moreover, we can assume that nR is an integer for n large enough.

- *Random codebook.* A *codebook* is a family \mathbf{b} of $|J| = 2^{nR}$ sequences $m^n(j)$ and $x^k(j)$ indexed

by $j \in J$. A *random codebook* is the draw of a codebook from the marginal i.i.d. probability distributions $(\lambda_m)^{\otimes n}$ and $\mathbb{P}(x)^{\otimes n}$. The selected codebook is known by the encoder and the decoder.

- *Encoding function.* The encoder observes the sequence of states $\omega^n \in \Omega^n$. It finds an index $j \in J$ such that the sequences $(\omega^n, m^n(j)) \in A_\delta$ are jointly typical, i.e. satisfy Equation (25). The encoder sends the sequence $x^k(j)$ corresponding to the index $j \in J$.
- *Decoding function.* The decoder observes the sequence of channel output $y^k \in Y^k$. It finds an index $\hat{j} \in J$ such that the sequences $(x^k(\hat{j}), y^k) \in A_\delta$ are jointly typical, i.e. satisfy Equation (26). The decoder decodes the sequence $m^n(\hat{j})$.
- *Error Event.* We introduce the indicator of error $E_\delta \in \{0, 1\}$ defined as follows:

$$E_\delta = \begin{cases} 0 & \text{if } j = \hat{j} \text{ and } (\omega^n, m^n, x^k, y^k) \in A_\delta, \\ 1 & \text{if } j \neq \hat{j} \text{ or } (\omega^n, m^n, x^k, y^k) \notin A_\delta. \end{cases} \quad (29)$$

An error $E_\delta = 1$ occurs in the coding process if: 1) the indices $j \in J$ and $\hat{j} \in J$ are not equal or 2) the sequences of symbols $(\omega^n, m^n, x^k, y^k) \notin A_\delta$, i.e. are not jointly typical.

An important result in information theory is that the expected probability of error over the random codebook is small.

Expected error probability. For all $\varepsilon_2 > 0$, for all $\eta > 0$, there exists a $\bar{\delta} > 0$, for all $\delta \leq \bar{\delta}$ there exists \bar{n}, \bar{k} such that for all $n \geq \bar{n}, k \geq \bar{k}$ and $|\frac{k}{n} - r| \leq \varepsilon_2$, the expected probability of the following error events are bounded by ε_2 :

$$\mathbb{E} \left[\mathbb{P}_b \left(\forall j \in J, (\omega^n, m^n(j)) \notin A_\delta \right) \right] \leq \varepsilon_2, \quad (30)$$

$$\mathbb{E} \left[\mathbb{P}_b \left(\exists j' \neq j, \text{ s.t. } (y^k, x^k(j')) \in A_\delta \right) \right] \leq \varepsilon_2. \quad (31)$$

- Equation (30) comes from Equation (27) and the Covering Lemma A.18, (Gamal and Kim, 2011, Lemma 3.3, p. 62).

- Equation (31) comes from Equation (28) and the Packing Lemma A.19, (Gamal and Kim, 2011, Lemma 3.1, p. 46).

If the expected probability of error is small over the codebooks, then it has to be small for

at least one codebook. Following a standard analysis of the error probability, (Gamal and Kim, 2011, pp. 42–43, 60–61), Equations (30), (31) imply that:

$$\forall \varepsilon_2 > 0, \forall \eta > 0, \exists \bar{\delta} > 0, \forall \delta \leq \bar{\delta}, \exists \bar{n}, \bar{k}, \forall n \geq \bar{n}, \forall k \geq \bar{k}, \left| \frac{k}{n} - r \right| \leq \varepsilon_2, \exists \mathbf{b}^*, \text{ s.t. } \mathbb{P}_{\mathbf{b}^*}(E_\delta = 1) \leq \varepsilon_2. \quad (32)$$

The strategy σ of the sender consists in using this codebook \mathbf{b}^* in order to find the sequence $m^n(j)$ which is jointly typical with ω^n , and in sending the sequence $x^k(j)$. By construction, this satisfies Equation (32), i.e. it has a low probability of error.

Part 2. Control of the Beliefs. The previous construction has the property that the decoder who uses the decoding schemes, makes an error with small probability. Now, the receiver needs not use the decoding scheme. Actually, the receiver calculates the posterior belief on the sequence of states ω^n , given y^k . Our contribution is to show that those beliefs are close to the prescribed beliefs ν_m at most stages. We have the following chain of inequalities:

$$\begin{aligned} & \mathbb{E}_\sigma \left[\frac{1}{n} \sum_{t=1}^n D\left(\nu_{t,y^k}^\sigma \parallel \nu_{m_t}\right) \mid E_\delta = 0 \right] \\ &= \sum_{m^n, y^k} \mathbb{P}_\sigma(m^n, y^k | E_\delta = 0) \cdot \frac{1}{n} \sum_{t=1}^n D\left(\nu_{t,y^k}^\sigma \parallel \nu_{m_t}\right) \end{aligned} \quad (33)$$

$$= \frac{1}{n} \sum_{(\omega^n, m^n, y^k) \in A_\delta} \mathbb{P}_\sigma(\omega^n, m^n, y^k | E_\delta = 0) \cdot \log_2 \frac{1}{\prod_{t=1}^n \nu_{m_t}(\omega_t)} - \frac{1}{n} \sum_{t=1}^n H(\omega_t | \mathbf{y}^k, E_\delta = 0) \quad (34)$$

$$\leq \frac{1}{n} \sum_{(\omega^n, m^n, y^k) \in A_\delta} \mathbb{P}_\sigma(\omega^n, m^n, y^k | E_\delta = 0) \cdot \log_2 \frac{1}{\prod_{t=1}^n \nu_{m_t}(\omega_t)} - \frac{1}{n} \sum_{t=1}^n H(\omega_t | \mathbf{m}^n, \mathbf{y}^k, E_\delta = 0) \quad (35)$$

$$\leq \frac{1}{n} \sum_{(\omega^n, m^n, y^k) \in A_\delta} \mathbb{P}_\sigma(\omega^n, m^n, y^k | E_\delta = 0) \cdot n \cdot \left(H(\omega | \mathbf{m}) + \delta \right) - \frac{1}{n} H(\omega^n | \mathbf{m}^n, \mathbf{y}^k, E_\delta = 0) \quad (36)$$

$$\leq \frac{1}{n} I(\omega^n; \mathbf{m}^n, \mathbf{y}^k | E_\delta = 0) - I(\omega; \mathbf{m}) + \delta + \frac{1}{n} + \log_2 |\Omega| \cdot \mathbb{P}_\sigma(E_\delta = 1) \quad (37)$$

$$\leq \frac{1}{n} I(\omega^n; \mathbf{m}^n | E_\delta = 0) - I(\omega; \mathbf{m}) + \delta + \frac{2}{n} + 2 \log_2 |\Omega| \cdot \mathbb{P}_\sigma(E_\delta = 1) \quad (38)$$

$$\leq \eta + \delta + \frac{2}{n} + 2 \log_2 |\Omega| \cdot \mathbb{P}_\sigma(E_\delta = 1). \quad (39)$$

- Equation (33) comes from the definition of the expected K-L divergence.

- Equation (34) comes from the conditioning by $E_\delta = 0$, since the support of $\mathbb{P}_\sigma(\omega^n, m^n, y^k | E_\delta = 0)$ is included in A_δ .

- Equation (35) comes from the property of the entropy $H(\omega_t|\mathbf{m}^n, \mathbf{y}^k, E_\delta = 0) \leq H(\omega_t|\mathbf{y}^k, E_\delta = 0)$.

- Equation (36) comes from the property of typical sequences $(\omega^n, \mathbf{m}^n) \in A_\delta$, stated in Lemma A.20 and in Gamal and Kim (2011, Property 1, pp. 26), and the chain rule for entropy:

$$H(\omega^n|\mathbf{m}^n, \mathbf{y}^k, E_\delta = 0) \leq \sum_{t=1}^n H(\omega_t|\mathbf{m}^n, \mathbf{y}^k, E_\delta = 0).$$

- Equation (37) comes from Lemma A.22 (see section A.5), which implies

$$\frac{1}{n}H(\omega^n|E_\delta = 0) - \frac{1}{n}H(\omega^n) + \frac{1}{n} + \log_2 |\Omega| \cdot \mathbb{P}_\sigma(E_\delta = 1) \geq 0.$$

Adding this expression to Equation (36) yields Equation (37).

- Equation (38) comes from Lemma A.22 (see section A.5) which implies that

$$I(\omega^n; \mathbf{y}^k|\mathbf{m}^n, E_\delta = 0) \leq I(\omega^n; \mathbf{y}^k|\mathbf{m}^n) + 1 + n \cdot \log_2 |\Omega| \cdot \mathbb{P}_\sigma(E_\delta = 1) = 1 + n \cdot \log_2 |\Omega| \cdot \mathbb{P}_\sigma(E_\delta = 1),$$

where $I(\omega^n; \mathbf{y}^k|\mathbf{m}^n) = 0$, from the Markov chain property of the triple $(\omega^n, \mathbf{m}^n, \mathbf{y}^k)$.

- Equation (39) comes from the cardinality of the codebook¹⁰:

$$I(\omega^n; \mathbf{m}^n|E_\delta = 0) \leq H(\mathbf{m}^n) \leq \log_2 |J| = n \cdot R = n \cdot (I(\omega; \mathbf{m}) + \eta).$$

Then, we have:

$$\begin{aligned} 1 - \mathbb{P}_\sigma(B_{\alpha, \gamma, \delta}) &:= \mathbb{P}_\sigma(B_{\alpha, \gamma, \delta}^c) \\ &= \mathbb{P}_\sigma(E_\delta = 1)\mathbb{P}_\sigma(B_{\alpha, \gamma, \delta}^c|E_\delta = 1) + \mathbb{P}_\sigma(E_\delta = 0)\mathbb{P}_\sigma(B_{\alpha, \gamma, \delta}^c|E_\delta = 0) \\ &\leq \mathbb{P}_\sigma(E_\delta = 1) + \mathbb{P}_\sigma(B_{\alpha, \gamma, \delta}^c|E_\delta = 0) \\ &\leq \varepsilon_2 + \mathbb{P}_\sigma(B_{\alpha, \gamma, \delta}^c|E_\delta = 0). \end{aligned} \tag{40}$$

¹⁰The last argument is inspired by Merhav and Shamai, 2007, Equation (23), for the problem of “Information Rates Subject to State Masking”.

Moreover:

$$\begin{aligned} & \mathbb{P}_\sigma(B_{\alpha,\gamma,\delta}^c | E_\delta = 0) \\ &= \sum_{m^n, y^k} \mathbb{P}_\sigma\left((m^n, y^k) \in B_{\alpha,\gamma,\delta}^c \mid E_\delta = 0\right) \end{aligned} \quad (41)$$

$$= \sum_{m^n, y^k} \mathbb{P}_\sigma\left((m^n, y^k) \text{ s.t. } \frac{|T_\alpha(m^n, y^k)|}{n} < 1 - \gamma \mid E_\delta = 0\right) \quad (42)$$

$$= \mathbb{P}_\sigma\left(\frac{\#}{n} \left\{t, \text{ s.t. } D\left(\nu_{t,y^k}^\sigma \parallel \nu_{m_t}\right) \leq \frac{\alpha^2}{2 \ln 2}\right\} < 1 - \gamma \mid E_\delta = 0\right) \quad (43)$$

$$= \mathbb{P}_\sigma\left(\frac{\#}{n} \left\{t, \text{ s.t. } D\left(\nu_{t,y^k}^\sigma \parallel \nu_{m_t}\right) > \frac{\alpha^2}{2 \ln 2}\right\} \geq \gamma \mid E_\delta = 0\right) \quad (44)$$

$$\leq \frac{2 \ln 2}{\alpha^2 \gamma} \cdot \mathbb{E}_\sigma \left[\frac{1}{n} \sum_{t=1}^n D\left(\nu_{t,y^k}^\sigma \parallel \nu_{m_t}\right) \right] \quad (45)$$

$$\leq \frac{2 \ln 2}{\alpha^2 \gamma} \cdot \left(\eta + \delta + \frac{2}{n} + 2 \log_2 |\Omega| \cdot \mathbb{P}_\sigma(E_\delta = 1) \right). \quad (46)$$

- Equations (41) to (44) are simple reformulations.

- Equation (45) comes from a use of Markov's inequality, detailed in Lemma A.21 (see section A.5).

- Equation (46) comes from equation (39).

Combining equations (32), (40), and (46) we obtain the following statement:

$$\forall \varepsilon_3 > 0, \forall \alpha > 0, \forall \gamma > 0, \exists \bar{\eta}, \forall \eta \leq \bar{\eta}, \exists \bar{\delta}, \forall \delta \leq \bar{\delta}, \exists \bar{n}, \bar{k}, \forall n \geq \bar{n}, \forall k \geq \bar{k}, \left| \frac{k}{n} - r \right| \leq \varepsilon_3, \exists \sigma,$$

such that:

$$\mathbb{P}_\sigma(B_{\alpha,\gamma,\delta}^c) \leq 2 \cdot \mathbb{P}_\sigma(E_\delta = 1) + \frac{2 \ln 2}{\alpha^2 \gamma} \cdot \left(\eta + \delta + \frac{2}{n} + 2 \log_2 |\Omega| \cdot \mathbb{P}_\sigma(E_\delta = 1) \right) \leq \varepsilon_3.$$

By choosing appropriately the “rate” $R \geq 0$ in (27) and (28) such as to make $\eta > 0$ small, we obtain the desired result:

$$\forall \varepsilon > 0, \forall \alpha > 0, \forall \gamma > 0, \exists \bar{\delta}, \forall \delta \leq \bar{\delta}, \exists \bar{n}, \bar{k}, \forall n \geq \bar{n}, \forall k \geq \bar{k}, \left| \frac{k}{n} - r \right| \leq \varepsilon, \exists \sigma,$$

such that $\mathbb{P}_\sigma(B_{\alpha,\gamma,\delta}^c) \leq \varepsilon$. □

A.5 Additional lemmas

The next three lemmas are standard results in information theory. They are recalled for the convenience of the reader.

Lemma A.18. (Covering lemma: compression of information source, Lemma 3.3, p. 62 in [Gamal and Kim, 2011](#))

Consider a random sequence ω^n with i.i.d. distribution $\mathbb{P}^{\otimes n}(\omega)$ and a family of 2^{nR} sequences $(m^n(j))_{j \in \{1, \dots, 2^{nR}\}}$ independently drawn from the i.i.d. distribution $\mathbb{P}^{\otimes n}(m)$. Assume that $R = I(\omega; \mathbf{m}) + \eta$ with $\eta > 0$.

For all $\varepsilon > 0$, there exists $\bar{\delta} > 0$, such that for all $\delta \leq \bar{\delta}$, there exists \bar{n} , such that for all $n \geq \bar{n}$:

$$\mathbb{P}\left(\forall j \in J, \quad (\omega^n, m^n(j)) \notin A_\delta\right) \leq \varepsilon.$$

Lemma A.19. (Packing lemma: transmission over a noisy channel, Lemma 3.1, p. 46 [Gamal and Kim, 2011](#))

Consider a random sequence y^k drawn with i.i.d. distribution $\mathbb{P}^{\otimes k}(y)$ and a family of 2^{kR} sequences $(x^k(j))_{j \in \{1, \dots, 2^{kR}\}}$ independently drawn from the i.i.d. distribution $\mathbb{P}^{\otimes k}(x)$. Assume that $R = I(\mathbf{x}; \mathbf{y}) - \eta$ with $\eta > 0$.

For all $\varepsilon > 0$, there exists $\bar{\delta} > 0$, such that for all $\delta \leq \bar{\delta}$, there exists \bar{k} , such that for all $k \geq \bar{k}$:

$$\mathbb{P}\left(\exists j \in J, \quad (x^k(j), y^k) \in A_\delta\right) \leq \varepsilon.$$

Lemma A.20 (Typical sequences, Property 1, p. 26 in [Gamal and Kim, 2011](#)). The typical sequences $(\omega^n, m^n) \in A_\delta$ satisfy:

$$\forall \delta_2 > 0, \exists \bar{\delta}_2 > 0, \forall \delta \leq \bar{\delta}_2, \forall n, \forall (\omega^n, m^n) \in A_\delta, \left| \frac{1}{n} \cdot \log_2 \frac{1}{\prod_{t=1}^n \mathbb{P}(\omega_t | m_t)} - H(\omega | \mathbf{m}) \right| \leq \delta_2,$$

where $\bar{\delta}_2 = \delta_2 \cdot H(\omega | \mathbf{m})$.

The next two lemmas are easy ancillary results that were used in the proofs and were omitted in the previous section to ease the reading.

Lemma A.21 (Markov's inequality). *For all $\varepsilon_1 > 0$, $\varepsilon_2 > 0$ we have:*

$$\mathbb{E}_\sigma \left[\frac{1}{n} \sum_{t=1}^n D \left(\mathbb{P}_\sigma(\boldsymbol{\omega}_t | \mathbf{y}^n, E_\delta = 0) \parallel \mathbb{P}(\boldsymbol{\omega}_t | \mathbf{m}_t) \right) \right] \leq \varepsilon_0 \quad (47)$$

$$\implies \mathbb{P}_{m^n, y^n} \left(\frac{\#}{n} \left\{ t, \text{ s.t. } D \left(\mathbb{P}_\sigma(\boldsymbol{\omega}_t | \mathbf{y}^n, E_\delta = 0) \parallel \mathbb{P}(\boldsymbol{\omega}_t | \mathbf{m}_t) \right) > \varepsilon_1 \right\} > \varepsilon_2 \right) \leq \frac{\varepsilon_0}{\varepsilon_1 \cdot \varepsilon_2}. \quad (48)$$

Proof. [Lemma A.21] We denote by $D_t = D(\mathbb{P}_\sigma(\boldsymbol{\omega}_t | \mathbf{y}^n, E_\delta = 0) \parallel \mathbb{P}(\boldsymbol{\omega}_t | \mathbf{m}_t))$ and $D^n = \{D_t\}_t$ the K-L divergence. We have that:

$$\mathbb{P} \left(\frac{\#}{n} \left\{ t, \text{ s.t. } D_t > \varepsilon_1 \right\} > \varepsilon_2 \right) = \mathbb{P} \left(\frac{1}{n} \cdot \sum_{t=1}^n \mathbf{1} \left\{ D_t > \varepsilon_1 \right\} > \varepsilon_2 \right) \quad (49)$$

$$\leq \frac{\mathbb{E} \left[\frac{1}{n} \cdot \sum_{t=1}^n \mathbf{1} \left\{ D_t > \varepsilon_1 \right\} \right]}{\varepsilon_2} \quad (50)$$

$$= \frac{\frac{1}{n} \cdot \sum_{t=1}^n \mathbb{E} \left[\mathbf{1} \left\{ D_t > \varepsilon_1 \right\} \right]}{\varepsilon_2} \quad (51)$$

$$= \frac{\frac{1}{n} \cdot \sum_{t=1}^n \mathbb{P} \left(D_t > \varepsilon_1 \right)}{\varepsilon_2} \quad (52)$$

$$\leq \frac{\frac{1}{n} \cdot \sum_{t=1}^n \frac{\mathbb{E}[D_t]}{\varepsilon_1}}{\varepsilon_2} \quad (53)$$

$$= \frac{1}{\varepsilon_1 \cdot \varepsilon_2} \cdot \mathbb{E} \left[\frac{1}{n} \cdot \sum_{t=1}^n D_t \right] \leq \frac{\varepsilon_0}{\varepsilon_1 \cdot \varepsilon_2}. \quad (54)$$

- Equations (49), (51), (52), (54) are reformulations of probabilities and expectations.

- Equations (50), (53), come from Markov's inequality $\mathbb{P}(X \geq \alpha) \leq \mathbb{E}[X]/\alpha$. \square

Lemma A.22. *Consider an i.i.d. random sequence $\boldsymbol{\omega}^n$. For all $\varepsilon > 0$, there exists $\bar{n} \in \mathbb{N}$ such that for all $n \geq \bar{n}$ we have:*

$$H(\boldsymbol{\omega}^n | E_\delta = 0) \geq n \cdot \left(H(\boldsymbol{\omega}) - \varepsilon \right). \quad (55)$$

Proof. [Lemma A.22]

$$H(\boldsymbol{\omega}^n | E_\delta = 0) = \frac{1}{\mathbb{P}(E_\delta = 0)} \cdot \left(H(\boldsymbol{\omega}^n | E_\delta = 1) - \mathbb{P}(E_\delta = 1) \cdot H(\boldsymbol{\omega}^n | E_\delta = 1) \right) \quad (56)$$

$$\geq H(\boldsymbol{\omega}^n | E_\delta) - \mathbb{P}(E_\delta = 1) \cdot H(\boldsymbol{\omega}^n | E_\delta = 1) \quad (57)$$

$$\geq H(\boldsymbol{\omega}^n) - H(E_\delta) - \mathbb{P}(E_\delta = 1) \cdot H(\boldsymbol{\omega}^n | E_\delta = 1) \quad (58)$$

$$\geq H(\boldsymbol{\omega}^n) - n \cdot \varepsilon. \quad (59)$$

- Equation (56) comes from the definition of the conditional entropy.
- Equation (57) comes from the property $\mathbb{P}(E_\delta = 0) \leq 1$.
- Equation (58) comes from the property $H(\boldsymbol{\omega}^n | E_\delta) = H(\boldsymbol{\omega}^n, E_\delta) - H(E_\delta) \geq H(\boldsymbol{\omega}^n) - H(E_\delta)$.
- Equation (59) comes from the i.i.d. property of the state ω and the definition of the error event $E_\delta = 1$. Hence, for all ε , there exists a $\bar{n} \in \mathbb{N}$ such that for all $n \geq \bar{n}$ we have: $H(\mathbb{P}(E_\delta = 1)) + \mathbb{P}(E_\delta = 1) \cdot \log_2 |\Omega| \leq \varepsilon$. □