Thanks to Ben Jones and Karl Floersch for clarifying discussions of the problem, and to Yan Zhang and Barnabé Monnot for helpful feedback on the presentation.

Dispute protocols are fundamentally about proving things: a proposition A

about events off chain needs to be decided on chain, so an adversarial game is set up on chain between a player affirming A and a player denying A

whose rules have been designed so that the correct player has a winning strategy.

Example.

Alice and Bob want a system to transfer an asset back and forth arbitrarily many times without any transactions on chain except for an initial deposit and a final withdrawal. To do this, they agree that the asset will be held on chain by a contract C

, and that transferring ownership of the asset off chain will be signified by sending the other party a signature of the current time. Ownership of the asset by (respectively) Alice and Bob is signified by the propositions

and for (say) Alice to withdraw the asset on chain, she must prove the proposition A

by winning the following game, which is adjudicated by the contract C

:

- · Alice is asked for a time t
- , and a signature by Bob of t
- . If Alice fails to provide these in reasonable time, then the game ends and Bob wins; otherwise:
  - . Bob is asked for a time t'
- , where t \leq t'
- , and a signature by Alice of t'
- . If Bob fails to provide these in reasonable time, then the game ends and Alice wins; otherwise:
  - The game ends and Bob wins.

Given the resemblance between the game and the proposition it decides, a natural question one might ask is: Instead of inventing and implementing such games and strategies de novo for each new dispute protocol, can we derive them automatically from the propositions and proofs they are fundamentally about?

This question was first asked and investigated by Ben Jones and Karl Floersch in their work on <u>fredicate contracts</u>" and later "<u>optimistic game semantics</u>". Inspired by <u>dialogical logic</u>, their work specifies a "universal adjudicator" that can interpret any proposition as a game.

This post proposes a different approach to the same question, based on<u>type theory</u> and <u>ludics</u>, which enables us to specify simultaneously a "universal adjudicator" and a corresponding "universal advocate" that can interpret any proof

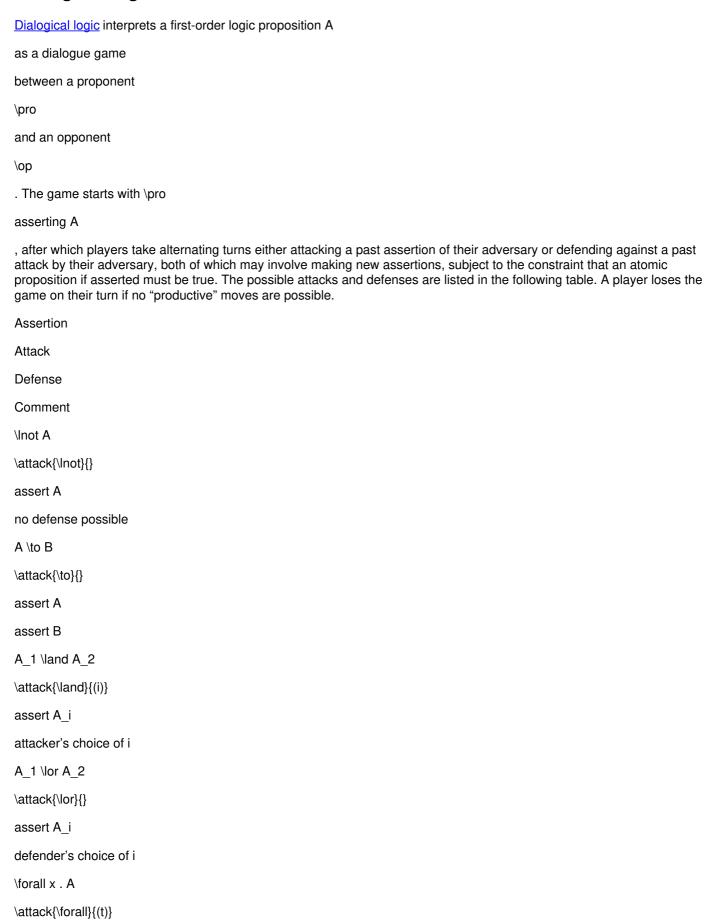
of a proposition as a winning strategy

for the corresponding game.

The first two sections are mostly review of background material (with some liberties taken to keep the presentation simple). Readers already familiar with propositions-as-types can skip directly to the third and last section titled "Games and strategies".

## **Dialogical logic**

assert A[x \mapsto t]



attacker's choice of t \exists x . A \attack{\exists}{} assert A[x \mapsto t] defender's choice of t Example. A play of the game corresponding to \lnot (A \land \lnot A) could go as follows (assuming the atomic proposition A is true; otherwise \op would lose the game on turn 4): 1. \pro starts, asserting \Inot (A \land \Inot A) 1. \op attacks move 1's assertion with \attack{\lnot}{} , asserting A \land \lnot A 1. \pro attacks move 2's assertion with \attack{\land}{(1)} 1. \op defends against move 3's attack, asserting A 1. \pro attacks move 2's assertion with \attack{\land}{(2)} 1. \op defends against move 5's attack, asserting \lnot A 1. \pro attacks move 6's assertion with \attack{\lnot}{} , asserting A 1. \op loses the game These games have the necessary property that there is a winning \pro strategy for any true proposition and a winning \op

## Type theory

strategy for any false proposition.

If we want to write proofs and interpret them as winning strategies, we need a formal proof system. Arguably the best choice of formal proof system today is to be found in <u>type theory</u>. Most state-of-the-art software for theorem proving (e.g. <u>Coq</u>, <u>Lean</u>, <u>Idris</u>) is based on it, and even programmers with no background in formal logic but with some experience in functional languages will find it familiar.

The <u>philosophy</u> of type theory in essence is that proving a proposition means constructing a mathematical object that makes its truth evident. For example, proving that two things are isomorphic means constructing an isomorphism between them. In general, every proposition is understood as specifying a type of mathematical object to construct (which of course is sometimes possible and sometimes not).

Proposition
Evidence
Set analogy
A \to B
procedure to transform evidence for A
into evidence for B
B ^ A
A \land B
both evidence for A
and evidence for B
A \times B
A \lor B
either evidence for A
or evidence for B
A + B
Ntop
trivial
1
Nbot
impossible
0
(\Inot A
is understood to be synonymous with A \to \bot
.)
A type theory (with the indefinite article) is a formal language for expressing these constructions. There are many different type theories, and there is no all-encompassing <u>definition</u> of what exactly it means to be a type theory (as with the concept of "space" in mathematics), but generally speaking a type theory looks something like the following.
There is a grammar of types
(here just enough for propositional logic):
lem:lem:lem:lem:lem:lem:lem:lem:lem:lem:
and a grammar of expressions
or terms
representing constructions (treated as <u>abstract binding trees</u> ):
lem:lem:lem:lem:lem:lem:lem:lem:lem:lem:
The typing relation
or typing judgment

"e

```
is of type A
", written
e : A,
states that the construction expressed by the term e
satisfies the specification expressed by the type A
 . A construction can be parameterized by a set of variables x_1, x_2, \ldots
respectively assumed to be of some types A 1, A 2, \ldots
 , so the typing judgment is generalized to "e
is of type A
in context Gamma = x_1 : A_1, x_2 : A_2, Idots
", written
\Gamma \vdash e : A.
Typing rules
define when a typing judgment (the conclusion
 , appearing below the line) can be derived from other typing judgments (the premises
 , appearing above the line).
Introduction rules
govern how to produce things of a type:
\frac{\Gamma, x : A \cdot B}{\Gamma \cdot A \cdot B} : B {\Gamma \cdot A \cdot B} \to B} 
\frac {} {\Gamma \vdash \unitI : \top} \top_\text{I}
\frac {\Gamma \vdash e : A} {\Gamma \vdash \eitherIL{e} : A \lor B} \lor_\text{I1} \qquad \frac {\Gamma \vdash e : B}
{\Gamma}_{e} : A \setminus B \setminus B \setminus A \setminus B
Elimination rules
govern how to consume things of a type:
\frac {\Gamma \vdash e_0 : A \to B \qquad \Gamma \vdash e : A} {\Gamma \vdash \funE{e_0}{e} : B} \to_\text{E}}
\frac{\Gamma \cdot Gamma \cdot G
\land_\text{E}
f(x) = 1 + 100  \frac {\Gamma \vdash e_0 : A \lor B \qquad \Gamma, x : A \vdash e_1 : C \qquad \Gamma, y : B \vdash e_2 : C} {\Gamma \vdash e_1 : C \qquad \Gamma, y : B \vdash e_2 : C} {\Gamma \vdash e_1 : C \qquad \Gamma, y : B \vdash e_2 : C}
\vdash = E\{e_0\}\{x\}\{e_1\}\{y\}\{e_2\} : C\} \cdot \{e_1\}\{y\}\{e_2\} : C\}
\frac {\Gamma \vdash e_0 : \bot} {\Gamma \vdash \voidE{e_0} : C} \bot_\text{E}
Finally, structural rules
govern general features like the use of variables:
\frac {} {\Gamma, x : A \vdash x : A} \, \text{hyp}
Example.
Here is a term proving \lnot(A \land \lnot A)
\funl{p}{(\pairE{p}{x}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\columnwidth}{f}{\normalfont{\colu
```

and here is its typing derivation (suppressing unused hypotheses in contexts to save space):

\cfrac { { \cfrac { {\cfrac { \cfrac { \cfrac { \cfrac { \cfrac { \cfrac { \cfrac { {\cfrac { \cfrac { \cfrac { \cfrac { \cfrac { \cfrac { \cfrac { {\cfrac { \cfrac { \cfrac { \cfrac { \cfrac { \cfrac { \cfrac { {\cfrac { \cfrac { \cfrac { \cfrac { \cfrac { \cfrac { \cfrac { {\cfrac { \cfrac { \cfrac { \cfrac { \cfrac { \cfrac { \cfrac { {\cfrac { \cfrac { \cfrac { \cfrac { \cfrac { \cfrac { \cfrac { {\cfrac { \cfrac { \cfrac { \cfrac { \cfrac { \cfrac { \cfrac { {\cfrac { \cfrac { \cfrac { \cfrac { \cfrac { \cfrac { \cfrac { {\cfrac { \cfrac { \cfrac { \cfrac { \cfrac { \cfrac { \cfrac { { 

## Computation

When different terms express the same construction is a <u>notoriously subtle question</u>, but at minimum there is a <u>congruence</u>,

```
historically named \beta
-reduction
, generated by "cancellations" of introductions and eliminations:
{y}{e_2} \ \unitE{\unitI}{e'} &\betaEq e' \ \eitherE{\eitherIL{e}}{x}{e_1'}{y}{e_2'} &\betaEq \subi{e_1'}{x}{e} \
\left(eitherE(eitherIR(e))(x)(e_1')(y)(e_2') \right) \le \left(e_2'(x)(e_2')(x)(e_2')(x)(e_2')(x)(e_2')(x)(e_2') \right)
where \subi{e'}{x}{e}
denotes the (capture-avoiding) substitution of e
for x
in e'
Simplifying terms with \beta
-reduction converts an implicit
representation of a construction (like "1 + 7 + 49 + 343") into an explicit
one (like "400"), a process which can be thought of philosophically as simulating the mental act of realizing the construction.
A (closed) term that is fully reduced (ignoring subterms under binders) is a value
\frac {} {\val{\funl{x}{e}}}}
\label{eq:conditional} $$ \operatorname{\langle v_1} \quad \operatorname{\langle v_2}} {\operatorname{\langle v_1}_{v_2}} $$
\frac {} {\val{\unitl}}}
\frac{\langle va|\{v\}} {\langle va|\{v\}\}} 
The next step of reduction is not unique in general, and is disambiguated with a reduction relation
or operational semantics
{\int F(v_0)e'}} \quad \frac{v} {\left(\int F(v_0)e'}} \quad \frac{v} {\left(\int F(v_0)e'} \right)}{\left(\int F(v_0)e'} \quad \left(\int F(v_0)e'} \quad \left(\int F(v_0)e'\right) \quad \left(\int F(v_0)e' \quad \left(\int F(v_0)e'
{\left[v_1}{e_2}\right]}
\frac {\step{e_0}{e_0'}} {\step{\pairE{e_0}{x}{y}{e}}}{\pairE{e_0'}{x}{y}{e}}}
\frac {\val{v_1} \quad \val{v_2}} {\step{\pairE{\pairI{v_1}{x}{y}{e}}}{\subii{e}{x}{v_1}{y}{e^2}}}
```

\frac {\step{e\_0}{e\_0'}} {\step{\unitE{e\_0}{e}}}{\unitE{e\_0'}{e}}}

\frac {} {\step{\unitE{\unitI}{e}}}{e}}

\frac {\val{v}} {\step{\eitherE{\eitherIL{v}}{x}{e\_1}{y}{e\_2}}{\subi{e\_1}{x}{v}}}

This process can be automated, effectively making a type theory a (terminating!) programming language. This is the celebrated <u>propositions-as-types/proofs-as-programs correspondence</u>.

## Games and strategies

With a formal proof system now in hand, we could attempt to transform proofs (i.e. terms) into strategies for the dialogue games described previously. This can indeed be done with more <u>machinery</u>, however I want to propose instead a simpler approach inspired by <u>ludics</u>, which takes advantage of the intrinsic computational aspect of type theory to interpret types as games that are naturally suited to interpreting terms as strategies.

To recapitulate, the goal is to transform:

· a type A

into a game between players \pro

and \op

• a term e \pro : A

into a winning strategy for \pro

a term e\_\op : A \to \bot

into a winning strategy for \op

A trivial noninteractive solution would be simply to require players to provide a term to win, but this runs into problems with atomic propositions about real events, for example "a hash preimage of \texttt{0x123}

has been revealed". Since evidence for this event would be the preimage in question, it naturally corresponds to a primitive type S

(for "secret") whose values are hash preimages of \texttt{0x123}

. The problem with the trivial game for this type is that while \pro

can win when it is true, \op

can never win even when it is false because there is no term of type S \to \bot

From a computational point of view, we can still ask if the behavior specified by this type can be "safely" implemented by "unsafe" code (like Rust's unsafe blocks). After all, as long as constructing a value of type S is really impossible, if there was a function \funl{x}{e}

```
of type S \to \bot
, then the body e
would be unreachable code anyway, so even ill-formed code could not cause a runtime error.
To allow such implementations, we introduce an "unsafe primitive" called a hole
(corresponding in ludics to the daimon
\maltese
):
\begin{alignat}{1} &\bbox[yellow]{\text{Holes} \; \hole{h}} \ &\text{Variables} \; x, y, z \ &\text{Expressions} \; e, v \; ::= \
&\gquad \; x &\bbox[yellow]{\sep \hole{h}} \ &\sep \funl{x}{e} &\sep \funE{e 0}{e} \ &\sep \pairl{e 1}{e 2} &\sep \pairE{e 0}{x}
{y}{e} \ &\sep \unitI &\sep \unitE{e_0}{e} \ &\sep \eitherIL{e} \sep \eitherIR{e} &\sep \eitherE{e_0}{x}{e_1}{y}{e_2} \ & &\sep
\voidE{e_0} \end{alignat}
A hole can occur at any type:
\frac {} {\Gamma \vdash \hole{h} : A} \, \text{hole}
and its behavior is to abort the evaluation (like a fatal exception):
\frac {} {\stuck{\hole{h}}{h}}
\frac {\stuck{e_1}{h}} {\stuck{\pairl{e_1}{e_2}}{h}} \qquad \frac {\val{v_1} \qquad \stuck{e_2}{h}} {\stuck{\pairl{v_1}{e_2}}{h}}
\frac {\stuck{e_0}{h}} {\stuck{\pairE{e_0}{x}{y}{e}}{h}}
\frac {\stuck{e_0}{h}} {\stuck{\unitE{e_0}{e}}}{h}}
\frac \\stuck{e}{h}} \\stuck{\eitherIL{e}}{h}} \\qquad \\frac \\stuck{e}{h}} \\stuck{\eitherIR{e}}{h}}
\frac {\stuck{e_0}{h}} {\stuck{\eitherE{e_0}{x}{e_1}{y}{e_2}}{h}}
Now terms eventually either reduce to a value or get stuck on a hole.
Theorem.
If e: A
then e \rightsquigarrow \ldots \rightsquigarrow e'
where e \sim e'
and e': A
and either e' \textsf{ value}
or e' \textsf{ stuck } \hole{h}
for some \hole{h}
contained in e
Let us say that a term not containing holes is total
, a term possibly containing holes is partial
, and a partial term e
is safe
when the holes it contains are unreachable by evaluation, i.e. there is no partial term e'
```

containing e

that gets stuck on a hole contained in e

- . In particular, total terms are trivially safe. In the case of the type S \to \bot
- , there is no total term, but if (and only if) no value of type S

has been revealed, then there is a safe partial term \funl{x}{\hole{h}}

.

The intention is that safe partial terms should also yield winning strategies. Now the goal is to transform:

• a type A

into a game between players \pro

and \op

· a safe partial

term e\_\pro : A

into a winning strategy for \pro

a safe partial

term e\_\op : A \to \bot

into a winning strategy for \op

This is accomplished by the following game for a type A

:

• \pro

and \op

respectively provide partial terms\begin{alignat}{1} e\_\pro &: A \ e\_\op &: A \ \to \bot \end{alignat}

• The partial term\funE{e\_\op}{e\_\pro}: \bot

is reduced until it gets stuck on a hole contained in either e\_\pro

or e\_\op

(which necessarily happens since there is no value of type \bot

to which it can reduce).

• The player on whose hole the reduction is stuck loses the game.