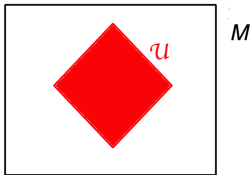


Lecture 3: A Background Independent Algebra in Quantum Gravity

Edward Witten

In ordinary quantum mechanics, we usually consider the observer to be outside the system that is being observed. This is problematical in the presence of gravity, most obviously in the case of a closed universe: No one can look at a closed universe from outside.

In ordinary quantum field theory, as discussed in the last few days, we can pick an arbitrary region \mathcal{U} in spacetime and define an algebra of observables in region \mathcal{U} :



There are problems with this in the presence of gravity: With spacetime fluctuating, it is in general hard to explain what we mean by the region \mathcal{U} . This only makes sense for particular regions that can be defined invariantly (for example, the exterior of a black hole horizon). But more fundamentally, why do we want to define an algebra unless it is the algebra of observables available to someone who lives in the spacetime?

So I want to construct an algebra that describes the measurements made by an observer. I will assume that the observer knows the laws of nature but has no knowledge of the state of the universe except whatever is gleaned from observation. (The second part fits our situation in the universe, but the first does not, since in the last few centuries we have been using our observations to learn the laws of nature as well as learning the state of the universe, i.e. part of what is usually called cosmology. I make the first assumption in part because it would be much harder to model an observer who is trying to learn the laws of nature. Also I suspect we should make this assumption if we want to arrive at the Bekenstein-Hawking entropy.)

The algebra will depend on the laws of nature, but it is required to be *universal* and *background independent*, meaning that it is defined once and for all without any knowledge of the specific spacetime in which the observer is living. Different spacetimes will correspond to different representations of the same algebra.

Some relevant papers:

Algebras of operators outside a black hole horizon

Leutheusser and Liu (2021)

EW (2021), Chandrasekharan, Penington, and EW (2022)

Algebra for a static patch in de Sitter space:

(*) Chandrasekharan, Longo, Penington, and EW (2022)

In JT gravity with negative cosmological constant

Penington and EW (2023), Kolchmeyer (2023)

In a general diamond-like region

Jensen, Sorce, and Speranza (2023)

Some recent developments

C. H. Chen and Penington (2024)

Kudler-Flam, Leutheusser, and Satishchandran (2024).

We expect that in a full theory of quantum gravity, an observer cannot be introduced from outside but must be described by the theory. What it means then to assume the presence of an observer is that we define an algebra that makes sense in a subspace of states in which an observer is present. We don't try to define an algebra that makes sense in all states.

First let us describe the situation in the absence of gravity. The observer propagates in a spacetime M on a geodesic γ :



The worldline is parametrized by proper time τ . As in classic work of Unruh (1976), the observer measures along γ , for example, a scalar field ϕ , or the electromagnetic field $F_{\mu\nu}$, or the Riemann tensor $R_{\mu\nu\alpha\beta}$, as well as their covariant derivatives in normal directions.

An elementary but not so well known point is that smearing a local field ϕ along a timelike curve is indeed sufficient to obtain a (densely defined, unbounded) operator. Hence there does rigorously exist an algebra of operators defined by smearing along the timelike worldline γ . By contrast, smearing ϕ in space is not effective at defining an operator, unless ϕ has rather low dimension. (As an important special case of this, it is not possible to define an operator by smearing a local field on a cut of a black hole or cosmological horizon.)

Let me take a moment to explain this. First the reason that smearing is necessary is that if $|\Omega\rangle$ is a Hilbert space state (for example the vacuum) and $\phi(x)$ is a local “operator,” then $\phi(x)|\Omega\rangle$ is never a Hilbert space state since its norm is

$$\langle\Omega|\phi^\dagger(x)\phi(x)|\Omega\rangle = \lim_{y\rightarrow x}\langle\Omega|\phi^\dagger(x)\phi(y)|\Omega\rangle = \infty.$$

Does smearing in d spatial dimensions help? Only if ϕ has rather low dimension: if

$$\phi_f = \int d^d x f(x)\phi(x)$$

is smeared in d spatial dimensions, then the norm of $\phi_f|\Omega\rangle$ is

$$\langle\Omega|\phi_f^\dagger\phi_f|\Omega\rangle = \int d^d x d^d y \bar{f}(x)f(y)\langle\Omega|\phi^\dagger(x)\phi(y)|\Omega\rangle.$$

If ϕ has dimension Δ as measured in the ultraviolet,

$$\phi^\dagger(x)\phi(y) \sim |x-y|^{-2\Delta} + \text{less singular terms}$$

then $\langle\Omega|\phi_f^\dagger\phi_f|\Omega\rangle < \infty$ if and only if

$$2\Delta < d.$$

For example, in QCD in the real world, $d = 3$, the minimum value of Δ is 3, so this is an example in which no operator can be defined by smearing in spatial directions.

A typical example of this came up yesterday. The quantity

$$K_R = \int_{x \geq 0} dx d\vec{y} x T_{00}(x, \vec{y})$$

is obtained by smearing in space only the local operator T_{00} . As T_{00} has dimension $\Delta = d + 1$, thus $2\Delta = 2d + 2 > d$, it should not be a surprise that K_R does not make sense as a Hilbert space operator.

Instead, smearing in time does make a true operator regardless of the value of Δ . The Feynman $i\epsilon$ is crucial: we have

$$\phi^\dagger(\vec{x}, t)\phi(\vec{x}, t') \sim (t - t' - i\epsilon)^{-2\Delta} + \dots$$

Now define

$$\phi_f(\vec{x}) = \int dt f(t)\phi(\vec{x}, t)$$

where f is a smooth smearing function of compact support. The norm of the state $\phi_f(\vec{x})|\Omega\rangle$ is then

$$\langle \phi_f(\vec{x})\Omega | \phi_f(\vec{x})\Omega \rangle = \int dt dt' \bar{f}(t) f(t') (t - t' - i\epsilon)^{-2\Delta}$$

(plus terms that are similar but less singular, from higher order terms in the OPE). This integral is clearly convergent for $\epsilon > 0$ and I claim that it has a finite limit for $\epsilon \rightarrow 0$.

For this, we just write

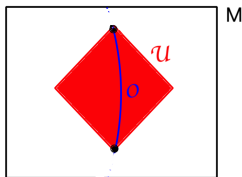
$$(t - t' - i\epsilon)^{-2\Delta} = C_n \frac{\partial^n}{\partial t^n} (t - t' - i\epsilon)^{n-2\Delta}$$

(for arbitrary integer $n > 0$ and a suitable constant C_n) so

$$\begin{aligned} \langle \phi_f(\vec{x}) \Omega | \phi_f(\vec{x}) \Omega \rangle &= C_n \int dt dt' \bar{f}(t) f(t') \frac{\partial^n}{\partial t^n} (t - t' - i\epsilon)^{n-2\Delta} \\ &= (-1)^n C_n \int dt dt' \frac{\partial^n \bar{f}(t)}{\partial t^n} f(t') (t - t' - i\epsilon)^{n-2\Delta}. \end{aligned}$$

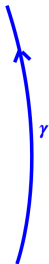
For sufficiently large n , this is obviously finite as $n \rightarrow 0$.

So it is possible to define an algebra of operators smeared along a timelike worldline γ . But what is this algebra? This question is answered by the “timelike tube” theorem (Borchers 1961; Araki 1963; Strohmaier 2000; Strohmaier and EW 2023), which is a close cousin of HKLL reconstruction in AdS/CFT duality. In quantum field theory without gravity, the algebra of operators along a timelike worldline γ is equivalent to the algebra of operators in a certain open set $\mathcal{E}(\gamma)$:



So in the absence of gravity, the algebra $\mathcal{A}(\gamma)$ of operators along a timelike geodesic γ is equivalent to the algebra of operators in a certain open set. Hence $\mathcal{A}(\gamma)$ is a reasonable substitute for the algebra of a region, which we usually consider in the absence of gravity, and appears to make more sense when gravity is included.

To understand more concretely what is $\mathcal{A}(\gamma)$ in the presence of gravity, let us focus on a particular observable, say $\phi(x(\tau))$ for a scalar field ϕ ; I will abbreviate this as $\phi(\tau)$. When we take gravity to be dynamical, we have to consider that the same worldline can be embedded in a given spacetime in different ways, differing by $\tau \rightarrow \tau + \text{constant}$:



So $\phi(\tau)$ isn't by itself a meaningful observable: we need to introduce the observer's degrees of freedom and define τ relative to the observer's clock.

In a minimal model, we equip the observer with a Hamiltonian $H_{\text{obs}} = mc^2 + q$, and a canonical variable $p = -i\frac{d}{dq}$. However, it turns out that it is better to assume that the observer energy is bounded below, say $q \geq 0$ (so m is the observer's rest mass). We then only allow operators that preserve this condition, so for example e^{-ip} , which does not preserve $q \geq 0$, should be replaced with $\Pi e^{-ip} \Pi$, where $\Pi = \Theta(q)$ is the projection operator onto $q \geq 0$.

We now want to allow only operators that commute with

$$\hat{H} = H_{\text{bulk}} + H_{\text{obs}},$$

where H_{bulk} is (any) gravitational constraint operator that generates a shift of τ along the worldline. An operator that commutes with \hat{H} is invariant under a spacetime diffeomorphism that moves the observer worldline forward in time, together with a time translation of the observer's system.

How do we find operators that commute with $\hat{H} = H + H_{\text{obs}} = H_{\text{bulk}} + H_{\text{obs}}$? Since

$$[H_{\text{bulk}}, \phi(\tau)] = -i\dot{\phi}(\tau),$$

we need

$$[q, \phi(\tau)] = i\dot{\phi}(\tau),$$

which we can achieve by just setting

$$\tau = p$$

or more generally

$$\tau = p + s$$

for a constant s .

So a typical allowed operator is $\phi(p + s)$, or more precisely

$$\hat{\phi}_s = \Pi \phi(p + s) \Pi = \Theta(q) \phi(p + s) \Theta(q).$$

In addition to these operators (with ϕ possibly replaced by any local field along the worldline such as the electromagnetic field or the Riemann tensor) there is one more obvious operator that commutes with \hat{H} , namely q itself. So we define an algebra \mathcal{A}_{obs} that is generated by the $\hat{\phi}_s$ as well as q .

The setup hopefully sounds “background independent,” since we described it without picking a background. However, background independence really depends on interpreting the formulas properly. We will not get background independence if we interpret $\hat{\phi}_S$ and q as Hilbert space operators. To get a Hilbert space on which $\hat{\phi}_S$ and q act, we have to pick a spacetime M in which the observer is propagating. Then we won't have background independence.

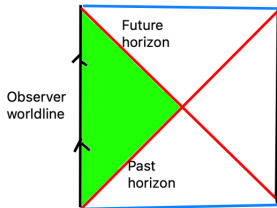
To get background independence, we have to think of \mathcal{A}_{obs} as an operator product algebra, rather than an algebra of Hilbert space operators. The algebras for different M 's are inequivalent representations of the same underlying operator product algebra.

In the absence of gravity, we would characterize the objects $\phi(\tau)$ by their universal short distance singularities:

$$\phi(\tau)\phi(\tau') \sim C(\tau - \tau' - i\epsilon)^{-2\Delta} + \dots .$$

This characterization does not require any knowledge about the quantum state. After coupling to gravity and including the observer and the constraint, the short distance expansion in powers of $\tau - \tau'$ becomes an expansion in $1/q$. We characterize \mathcal{A}_{obs} purely by the universal short distance or $1/q$ expansion of operator products. With that understanding, \mathcal{A}_{obs} is background-independent.

There is a very special case that turns out to be important. This is the case that M is an empty de Sitter space, with some positive value of the effective cosmological constant.



The green region is called a static patch, because it is invariant under a particular de Sitter generator H that advances the proper time of the observer.

In the absence of gravity, there is a distinguished de Sitter invariant state Ψ_{dS} such that correlation functions in this state are thermal at the de Sitter temperature $T_{\text{dS}} = 1/\beta_{\text{dS}}$ (Gibbons and Hawking; Figari, Nappi, and Hoegh-Krohn). For example, this means that two point functions $\langle \Psi_{\text{dS}} | \phi(\tau) \phi'(\tau') | \Psi_{\text{dS}} \rangle$ have two key properties:

(1) Time translation symmetry:

$$\langle \Psi_{\text{dS}} | \phi(\tau + s) \phi'(\tau' + s) | \Psi_{\text{dS}} \rangle = \langle \Psi_{\text{dS}} | \phi(\tau) \phi'(\tau') | \Psi_{\text{dS}} \rangle.$$

(2) The KMS condition, which says roughly:

$$\langle \Psi_{\text{dS}} | \phi(\tau) \phi'(0) | \Psi_{\text{dS}} \rangle = \langle \Psi_{\text{dS}} | \phi'(0) \phi(\tau - i\beta) | \Psi_{\text{dS}} \rangle.$$

(A precise statement involves holomorphy of the correlation function in a strip in the complex plane.)

Including gravity and the observer, we define a special state in which the observer energy has a thermal distribution at the de Sitter temperature

$$\Psi_{\max} = \Psi_{\text{dS}} e^{-\beta_{\text{dS}} q/2} \sqrt{\beta_{\text{dS}}},$$

and we replace operators $\phi(\tau)$ by “gravitationally dressed” operators $\hat{\phi}_s = \Pi \phi(p + s) \Pi$. Then a straightforward computation shows that

(1') We still have time-translation symmetry

$$\langle \Psi_{\max} | \hat{\phi}_s \hat{\phi}'_{s'} | \Psi_{\max} \rangle = \langle \Psi_{\max} | \hat{\phi}_{s+c} \hat{\phi}'_{s'+c} | \Psi_{\max} \rangle, \quad c \in \mathbb{R}.$$

(2') The KMS condition simplifies:

$$\langle \Psi_{\max} | \hat{\phi}_s \hat{\phi}'_{s'} | \Psi_{\max} \rangle = \langle \Psi_{\max} | \hat{\phi}'_{s'} \hat{\phi}_s | \Psi_{\max} \rangle.$$

Condition (2') tells us that if, for any $\mathbf{a} \in \mathcal{A}_{\text{obs}}$, we define

$$\text{Tr } \mathbf{a} = \langle \Psi_{\text{max}} | \mathbf{a} | \Psi_{\text{max}} \rangle,$$

then the function Tr does have the algebraic property of a trace:

$$\text{Tr } \mathbf{ab} = \text{Tr } \mathbf{ba}, \quad \mathbf{a}, \mathbf{b} \in \mathcal{A}_{\text{obs}}.$$

This function has the property that $\text{Tr } \mathbf{a}^\dagger \mathbf{a} > 0$ for all $\mathbf{a} \neq 0$, meaning in particular that it is “nondegenerate.” Note that if Ψ_{max} is normalized then

$$\text{Tr } 1 = 1.$$

Let \mathcal{H}_{dS} be the Hilbert space that we get by quantizing fields in de Sitter space (in perturbation theory). Let Ψ be any state in \mathcal{H}_{dS} and consider the function $\mathbf{a} \rightarrow \langle \Psi | \mathbf{a} | \Psi \rangle$, $\mathbf{a} \in \mathcal{A}_{\text{obs}}$. Roughly speaking, because \mathcal{A}_{obs} has the nondegenerate trace Tr , we can hope that there is a “density matrix” $\rho \in \mathcal{A}_{\text{obs}}$ such that

$$\langle \Psi | \mathbf{a} | \Psi \rangle = \text{Tr } \mathbf{a} \rho, \quad \mathbf{a} \in \mathcal{A}_{\text{obs}}.$$

Rather as in ordinary quantum mechanics, we expect ρ to be a positive element $\rho \in \mathcal{A}_{\text{obs}}$ with $\text{Tr } \rho = 1$. For example, let us find the density matrix of the state Ψ_{max} .

The definition of the trace makes it clear that the density matrix of the state Ψ_{\max} is $\sigma_{\max} = 1$, since to satisfy

$$\langle \Psi_{\max} | \mathbf{a} | \Psi_{\max} \rangle = \text{Tr } \mathbf{a} \sigma_{\max} \equiv \langle \Psi_{\max} | \mathbf{a} \sigma_{\max} | \Psi_{\max} \rangle,$$

we set

$$\sigma_{\max} = 1.$$

This means that Ψ_{\max} is “maximally mixed,” similar to a maximally mixed state in ordinary quantum mechanics whose density matrix is a multiple of the identity.

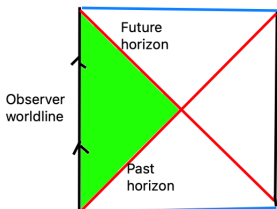
Now if $\mathbf{a} \in \mathcal{A}_{\text{obs}}$ is any operator, consider the state $\Psi_{\mathbf{a}} = \mathbf{a}\Psi_{\text{max}}$. It has a density matrix $\rho_{\Psi_{\mathbf{a}}} = \mathbf{a}\mathbf{a}^\dagger$, since for any $\mathbf{b} \in \mathcal{A}_{\text{obs}}$,

$$\langle \Psi_{\mathbf{a}} | \mathbf{b} | \Psi_{\mathbf{a}} \rangle = \langle \Psi_{\text{max}} | \mathbf{a}^\dagger \mathbf{b} \mathbf{a} | \Psi_{\text{max}} \rangle = \text{Tr } \mathbf{a}^\dagger \mathbf{b} \mathbf{a} = \text{Tr } \mathbf{b} \mathbf{a} \mathbf{a}^\dagger.$$

But states $\Psi_{\mathbf{a}}$ are dense in \mathcal{H}_{dS} – roughly by the Reeh-Schlieder theorem, which is the fundamental result about entanglement in quantum field theory. So a dense set of states have density matrices.

If we want *all* states to have density matrices, we need to take a useful further step. The Hilbert space \mathcal{H}_{dS} is the *closure* of a dense set of states $\mathbf{a}\Psi_{\text{dS}}$, so if we want *every* state in \mathcal{H}_{dS} to have a density matrix, we have to similarly take a closure of \mathcal{A}_{obs} . This closure, which is no longer background independent, can be defined as the von Neumann algebra generated by bounded operators in \mathcal{A}_{obs} . I will call the closure $\mathcal{A}_{\text{obs,dS}}$. Every state in \mathcal{H}_{dS} has a density matrix in (or technically, in general affiliated to) $\mathcal{A}_{\text{obs,dS}}$. It is in this step that von Neumann algebras enter the picture.

Now let us recall that Bekenstein and Hawking discovered that one should attribute an entropy to a black hole horizon. Not too long after, Gibbons and Hawking suggested that one should also attribute an entropy to a cosmological horizon such as the horizon experienced by an observer in de Sitter space:



However, a microscopic interpretation of what is meant by the entropy of a cosmological horizon has been obscure – even compared to the still largely mysterious black hole entropy. We can now give at least a partial answer to this question, at least for the case of de Sitter space.

Indeed once we know that every state of the observer algebra has a density matrix, we can define entropies as well. The von Neumann entropy of a density matrix ρ is as usual

$$S(\rho) = -\text{Tr } \rho \log \rho.$$

In ordinary quantum mechanics, a maximally mixed state has a density matrix that is a multiple of the identity, and it has the maximum possible von Neumann entropy. The analog here is Ψ_{max} , with density matrix $\sigma_{\text{max}} = 1$. It is clear that

$$S(\sigma_{\text{max}}) = -\text{Tr } 1 \log 1 = 0,$$

and by imitating an argument that in ordinary quantum mechanics proves that a maximally mixed state has maximum possible entropy, one can prove that every other density matrix $\rho \neq 1$ has strictly smaller entropy:

$$S(\rho) < 0.$$

One way to make this proof is as follows. Let $\rho \neq 1$ be any density matrix other than the identity. Then for $0 \leq t \leq 1$, $\rho_t = (1-t) + t\rho$ is also a density matrix. Let $f(t) = S(\rho_t) = -\text{Tr } \rho_t \log \rho_t$, so $S(\rho) = f(1)$. Then

$$f'(0) = 0, \quad f''(t) < 0 \quad \text{for } 0 \leq t \leq 1.$$

The first statement is almost immediate, and to prove the second, one uses $\log M = \int_0^\infty ds \left(\frac{1}{s} - \frac{1}{s+M} \right)$, which leads to

$$f''(t) = - \int_0^\infty ds \text{Tr } \frac{1}{s + \rho_t} (1-\rho) \frac{1}{s + \rho_t} (1-\rho) = - \int_0^\infty ds \text{Tr } B^2 < 0,$$

where B is the self-adjoint operator

$B = \left(\frac{1}{s+\rho_t} \right)^{1/2} (1-\rho) \left(\frac{1}{s+\rho_t} \right)^{1/2}$. Since $f'(0) = 0$, $f''(t) < 0$, we get $f(1) < 0$ so

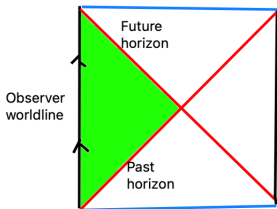
$$S(\rho) < 0.$$

Thus, the system consisting of an observer in a static patch in de Sitter space has a state of maximum entropy

$$\Psi_{\text{max}} = \Psi_{\text{dS}} e^{-\beta_{\text{dS}} q/2} \sqrt{\beta_{\text{dS}}},$$

consisting of empty de Sitter space with a thermal distribution of the observer energy. Why did this happen?

The original argument that empty de Sitter space has maximum entropy is due to Bousso (2000), who argued that this must be true, based on the Second Law of Thermodynamics, because the static patch is empty in the far future:



In the present context, we've defined the static patch by the presence of the observer, so by definition the observer doesn't leave the static patch even in the far future. But we can expect that in the far future the static patch will be empty except for the presence of the observer, and that the observer will be in thermal equilibrium with the bulk quantum fields, and that is what we see in the state Ψ_{max} . So the maximum entropy state that we found is the one suggested by Bousso's argument.

In a more general spacetime, I do not know an equally explicit definition of the entropy of a state of the observer algebra. However, there is a fairly reasonable conjecture, which is inspired by A. Wall's proof of the generalized second law (2011) and has some support from recent work of Chen and Penington (and still more recently Blommaert, Kudler-Flam, and Urbach). The idea is to interpret the Hartle-Hawking no boundary state as a sort of universal state of maximum entropy (generalizing empty de Sitter space, which has maximum entropy among states in a particular de Sitter spacetime). The expectation value in the no boundary state $\mathbf{a} \rightarrow \langle \Psi_{\text{HH}} | \mathbf{a} | \Psi_{\text{HH}} \rangle$ is a state of the observer algebra that I will denote as $\sigma(\mathbf{a})$. Then if $\mathbf{a} \rightarrow \rho(\mathbf{a})$ is any state of the observer algebra, then I suggest that the relative entropy between ρ and σ for the observer algebra gives - up to sign - a definition of the entropy of the state seen by the observer:

$$S(\rho) = -S(\rho|\sigma).$$

Here we should *not* try to normalize the no boundary state, rather we should just take it as it appears from the gravitational path integral. One reason for that is that it gives the “right answer.” Philosophically, another reason is that an observer living in a particular spacetime doesn’t know how to normalize the no boundary state, but may understand the no boundary state in the observer’s own spacetime.

The main evidence for the proposal is that in the case of a de Sitter like spacetime, it gives correctly the $A/4G$ term in the entropy. This statement reflects the $\log \sigma$ term in the relative entropy

$$-S(\rho|\sigma) = -\text{Tr}(\rho \log \rho - \rho \log \sigma) = S(\rho) + \text{Tr} \rho \log \sigma$$

and is a reinterpretation of the original calculation of Gibbons and Hawking. Beyond this, and apart from the analogy with the work of A. Wall, the main appeal of the proposal is that it makes sense universally as a definition of entropy in a cosmological setting.

The proposal only makes sense if it is true that the no boundary state is “tracial” for the observer algebra. This is certainly far from clear, but some arguments to that effect were given in a recent paper by Blommaert, Kudler-Flam, and Urbach 2505.14771.

The reason that we want the no boundary state to be tracial is that otherwise we cannot expect $-S(\rho|\sigma)$ to be a simple entropy measure:

$$-S(\rho|\sigma) = S(\rho) + \text{Tr } \rho \log \sigma,$$

and unless σ is tracial (so that the last term is a constant in any given spacetime) the last term spoils the entropic interpretation. It is the last term that for de Sitter spacetimes contributes the $A/4G$ term that the algebraic definition of entropy of the algebra doesn't "know" about.

A remark here is that it was an idealization to assume that the observer lives forever, and it was also an idealization to assume that the observer can make arbitrary measurements. I suspect that we need to make those idealizations if we want to get the Bekenstein-Gibbons-Hawking entropy as the entropy of a state of the observer algebra. As “entropy” measures information that could be obtained in principle about the state, but hasn’t been obtained, to the extent that we make the model of the observer more realistic, we will make the entropy smaller.

A lot of things are missing, for example:

- * What about an observer (or a civilization) that did not always exist?
- * In this presentation, I used a field theory language; can the discussion be generalized to string/M-theory?
- * Though the definitions make sense regardless, I want to remark that the presentation that I've given seems most natural for an observer who because of black hole or cosmological horizons cannot see the whole universe.
- * Finally is it possible to justify the proposal that I stated at the end?